

HC70A Winter 2006
Professor Bob Goldberg

Lecture #7
The Human Genome, Human Gene Diversity,
& Are there Races?

Themes/Concepts

- 1 Human Genomes
- 2 Mitochondrial Genome & Inheritance
- 3 Human Nuclear Genome
- 4 SNPs & Origins of Human Gene Diversity/Variation
- 5 Uses of SNPs
- 6 Haplotypes / Uses
- 7 Using Haplotypes to Trace Human Origins
- 8 Are there Races?
- 9 Skin Color Genes
- 10 Genetic Variation within & between Populations
- 11 "Where" does Most Human Gene Variation Reside?
- 12 Origins of Race Concept

— Sat 3/9/06

1hr only

HUMAN GENES ARE PRESENT
in TWO compartments --
The Nucleus & The Mitochondria

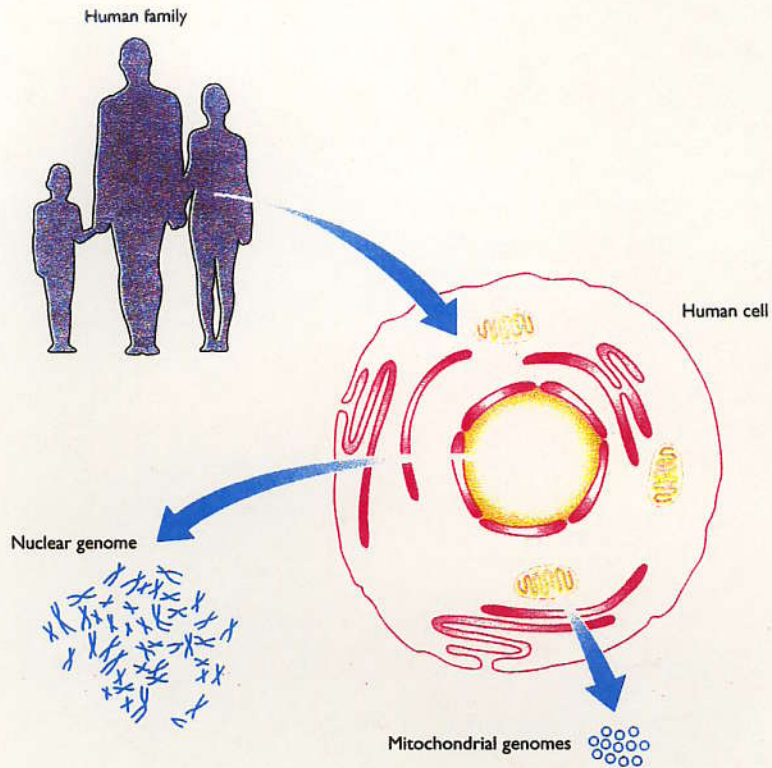


Figure 1.3 The nuclear and mitochondrial components of the human genome.

For more details on the anatomy of the human genome, see Section 6.1.

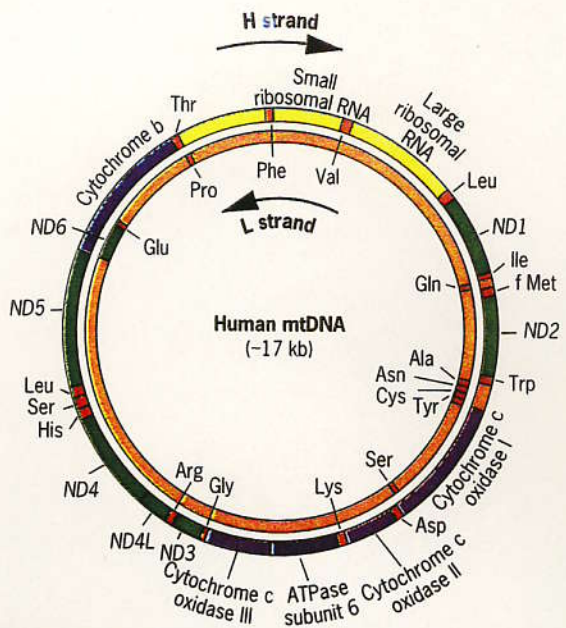
Genes in BOTH compartments are
critical for human development -

Features of the Human Nuclear and Mitochondrial Genomes

Table 9.1: The human nuclear and mitochondrial genomes

	Nuclear genome	Mitochondrial genome
Size	3200 Mb	16.6 kb
No. of different DNA molecules	23 (in XX cells) or 24 (in XY cells); all linear	One circular DNA molecule
Total no. of DNA molecules per cell	46 in diploid cells, but varies according to ploidy	Often several thousands (but variable – see Box 9.1)
Associated protein	Several classes of histone and nonhistone protein	Largely free of protein
No. of genes	~ 30 000–35 000	37
Gene density	~ 1/100 kb	1/0.45 kb
Repetitive DNA	Over 50% of genome, see Figure 9.1	Very little
Transcription	The great bulk of genes are transcribed individually (<i>monocistronic transcription units</i>)	Co-transcription of multiple genes from both the heavy and the light strands (<i>polycistronic transcription units</i>)
Introns	Found in most genes	Absent
% of coding DNA	~ 1.5%	~ 93%
Codon usage	See Figure 1.22	See Figure 1.22
Recombination	At least once for each pair of homologs at meiosis	Not evident
Inheritance	<u>Mendelian</u> for sequences on X and autosomes; paternal for sequences on Y	Exclusively <u>maternal</u>

The Mitochondrial Genome is A Small Circle containing only 37 Genes



MtDNA is a CIRCLE

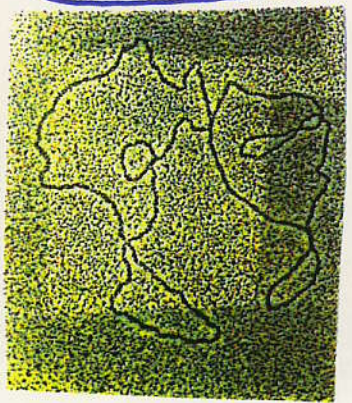
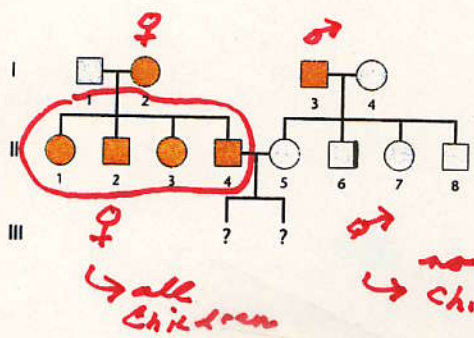


Figure 19.14 Map of human mtDNA showing the pattern of transcription. Genes on the inner circle are transcribed from the L strand of the DNA, whereas genes on the outer circle are transcribed from the H strand of the DNA. Arrows show the direction of transcription. ND1-6 are genes encoding subunits of the enzyme NADH reductase; the tRNA genes in the mtDNA are indicated by abbreviations for the amino acids.

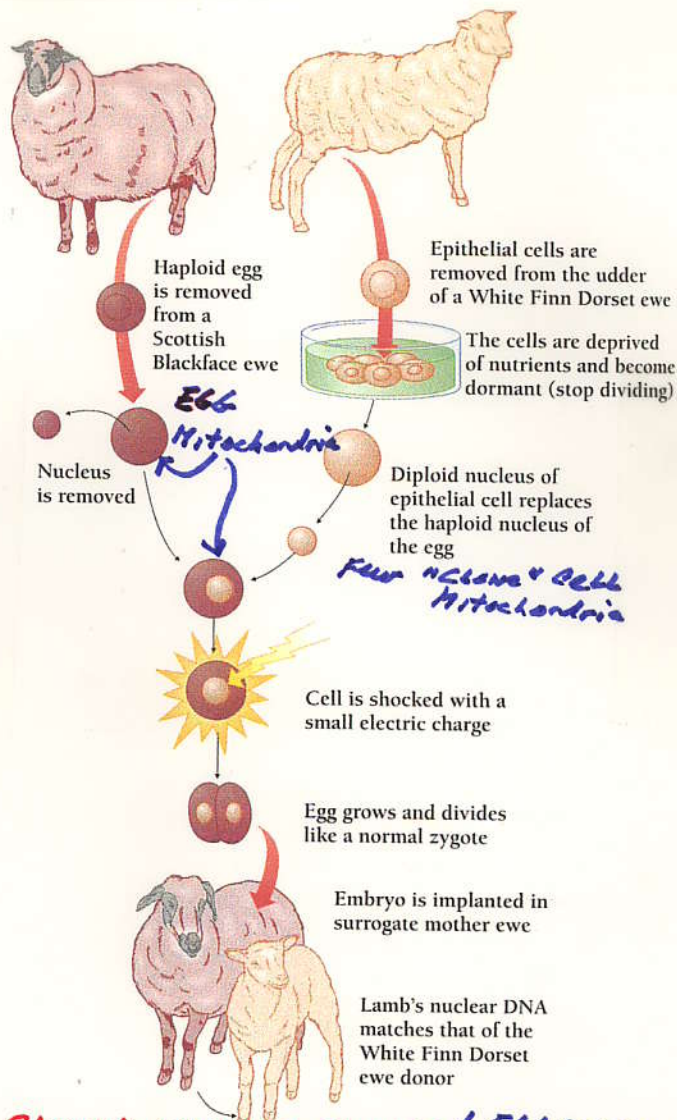
Mitochondrial Genes ARE Inherited MATERNALLY



PASSED DIRECTLY FROM MOTHER TO CHILDREN

Hypothesis to Explain?

IN A CLONING EXPERIMENT MOST OF MITOCHONDRIA COMES FROM EGG DONOR!



∴ CLONE NOT TECHNICALLY A CLONE WITH RESPECT TO MITOCHONDRIAL GENES

WILL BE MOSAIC - MOST MT FROM EGG DONOR - FEW MT FROM CELL USED FOR CLONE'S GENOME!!

CLONE'S NUCLEAR GENOME / Egg Donor Mt Genome!

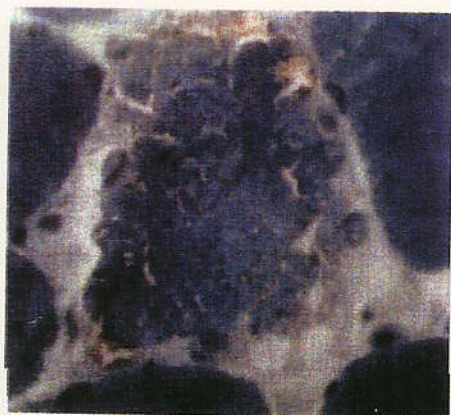
Figure 44-16 Cloning Dolly. The trick to cloning Dolly was to make differentiated cells less differentiated. By depriving the cultured udder cells of nutrients, the researchers induced the nuclei to enter a dormant state.

Another Potential Problem Source for Clone development!

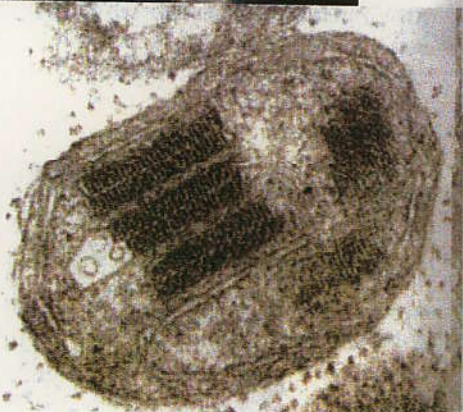
Several Mitochondrial Diseases Occur in Humans

In order for a human disorder to be attributable to genetically altered mitochondria, several criteria must be met.

1. Inheritance must exhibit a maternal rather than a Mendelian pattern.
2. The disorder must reflect a deficiency in the bioenergetic function of the organelle.
3. There must be a specific genetic mutation in one of the mitochondrial genes.



(a)



(b)

FIGURE 8.5 Mitochondrial myopathy in skeletal muscle cells of a patient with MERRF. Part (a) shows a ragged red fiber with abnormal mitochondria. Part (b) shows an abnormal mitochondrion revealing paracrystalline arrays within it.

Thus far, several cases are known to demonstrate these characteristics. For example, myoclonic epilepsy and ragged red fiber disease (MERRF) demonstrates a pattern of inheritance consistent with maternal inheritance. Only offspring of affected mothers inherit the disorder; the offspring of affected fathers are all normal. Individuals with this rare disorder express deafness, ataxia, and seizures. Both muscle fibers and mitochondria are affected; the aberrant mitochondria characterize what are described as ragged red fibers (RRFs) of skeletal muscle (Figure 8.5). Analysis of mtDNA has revealed a mutation in one of the mitochondrial genes encoding a transfer RNA. This genetic alteration apparently interferes with translation within the organelle, which in turn leads to the various manifestations of the disorder.

A second disorder, Leber's hereditary optic neuropathy (LHON), also exhibits maternal inheritance as well as mtDNA lesions. The disorder is characterized by sudden bilateral blindness. The average age of vision loss is 27, but onset is quite variable. Four mutations have been identified, all of which disrupt normal oxidative phosphorylation. Over 50 percent of cases are due to a mutation at a specific position in the mitochondrial gene encoding a subunit of NADH dehydrogenase so that the amino acid arginine is converted to histidine. This mutation is transmitted to all maternal offspring. It is interesting to note that in many instances of LHON, there is no family history; a significant number of cases appear to result from "new" mutations.

Individuals severely affected by a third disorder, Kearns-Sayre syndrome (KSS), lose their vision, undergo hearing loss, and display heart conditions. The genetic basis of KSS involves deletions at various positions within mtDNA. Many KSS patients are symptom-free as children but display progressive symptoms as adults. The proportion of mtDNAs that reveal deletions increases as the severity of symptoms increases.

The study of hereditary mitochondrial-based disorders provides insights into the importance and genetic basis of this organelle during normal development, as well as the relationship between mitochondrial function and neuromuscular disorders.

Such study has also suggested a hypothesis for aging based on the progressive accumulation of mtDNA mutations and the accompanying loss of mitochondrial function.

Mitochondrial Mutations
LEADING TO DISEASES

Table 16.1 Phenotypes associated with some mitochondrial mutations

Nucleotide changed	Mitochondrial component affected	Phenotype ^a
3460	ND1 of Complex I ^b	LHON
11778	ND4 of Complex I	LHON
14484	ND6 of Complex I	LHON
8993	ATP6 of Complex V ^b	NARP
3243	tRNA ^{Leu(UUR)} ^c	MELAS, PEO
3271	tRNA ^{Leu(UUR)}	MELAS
3291	tRNA ^{Leu(UUR)}	MELAS
3251	tRNA ^{Leu(UUR)}	PEO
3256	tRNA ^{Leu(UUR)}	PEO
5692	tRNA ^{Asn}	PEO
5703	tRNA ^{Asn}	PEO, myopathy
5814	tRNA ^{Cys}	Encephalopathy
8344	tRNA ^{Lys}	MERRF
8356	tRNA ^{Lys}	MERRF
9997	tRNA ^{Gly}	Cardiomyopathy
10006	tRNA ^{Gly}	PEO
12246	tRNA ^{Ser(AGY)} ^c	PEO
14709	tRNA ^{Glu}	Myopathy
15923	tRNA ^{Thr}	Fatal infantile multisystem disorder
15990	tRNA ^{Pro}	Myopathy

^aLHON Leber's hereditary optic neuropathy; NARP Neurogenic muscle weakness, ataxia, retinitis pigmentosa; MERRF Myoclonic epilepsy and ragged-red fiber syndrome; MELAS Mitochondrial myopathy, encephalopathy, lactic acidosis, stroke-like episodes; PEO Progressive external ophthalmoplegia

^bComplex I is NADH dehydrogenase. Complex V is ATP synthase.

^cIn tRNA^{Leu(UUR)}, the R stands for either A or G; in tRNA^{Ser(AGY)}, the Y stands for either T or C.

www.mitomap.org



MAP OF MITOCHONDRIAL GENES AND CORRESPONDING DISEASES

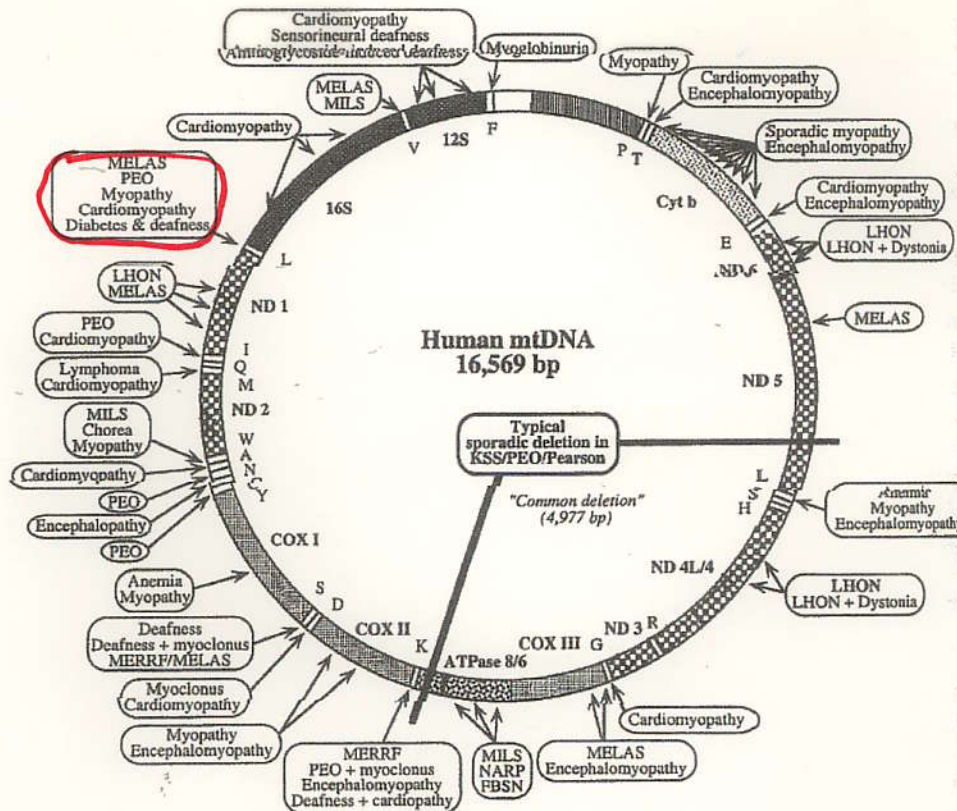


FIGURE 18.16. Morbidity map of the human mitochondrial genome. Abbreviations are for the genes encoding seven subunits of complex I (ND), three subunits of cytochrome c oxidase (COX), cytochrome b (Cyt b), and the two subunits of ATP synthase (ATPase 6 and 8). 12S and 16S refer to ribosomal RNAs; 22 transfer RNAs are identified by the one-letter codes for the corresponding amino acids. FBSN, familial bilateral striatal necrosis; KSS, Kearns-Sayre syndrome; LHON, Leber hereditary optic neuropathy; MELAS, mitochondrial encephalomyopathy, lactic acidosis, and strokelike episodes; MERRF, myoclonic epilepsy with ragged-red fibers; MILS, maternally inherited Leigh syndrome; NARP, neuropathy, ataxia, retinitis pigmentosa; PEO, progressive external ophthalmoplegia. From DiMauro and Schon (2001). Used with permission.

Mitochondrial Diseases ARE Inherited maternally

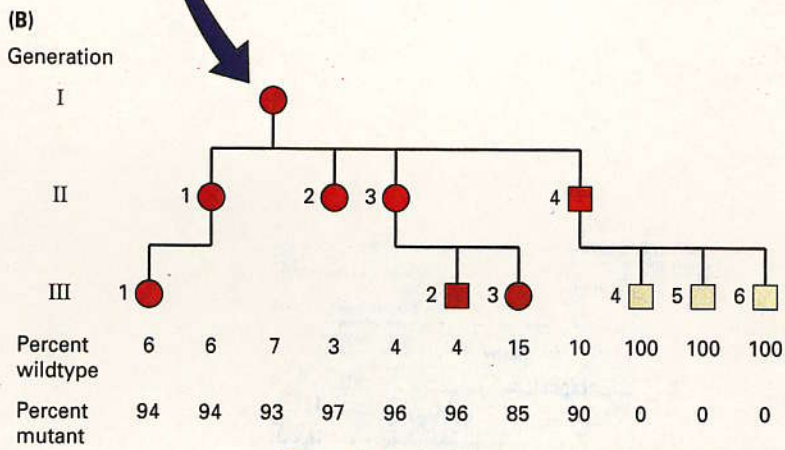
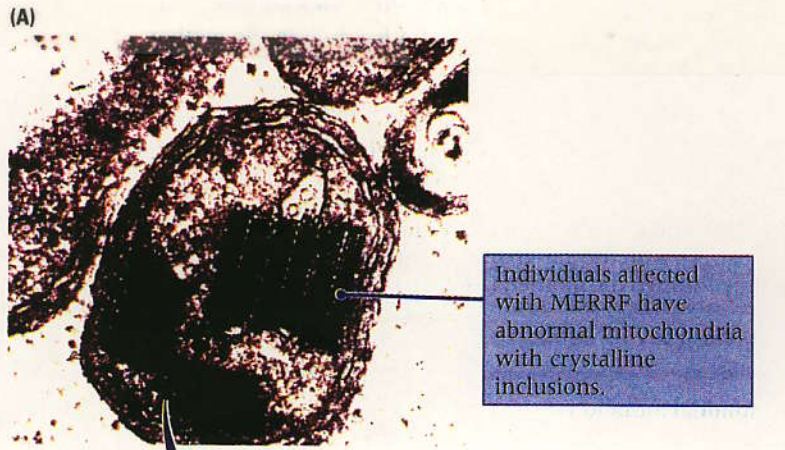


Figure 16.2 Inheritance of myoclonic epilepsy with ragged-red fiber disease (MERRF) in humans. (A) Electron micrograph of an abnormal MERRF mitochondrion containing paracrystalline inclusions. (B) The pedigree shows inheritance of MERRF in one family and the percentage of the mitochondria in each person found to be wildtype or mutant. [Micrograph courtesy of D. C. Wallace, from J. M. Shoffner, M. T. Lott, A. M. S. Lezza, P. Seibel, S. W. Ballinger, and D. C. Wallace. 1990. *Cell* 61: 931.]

Never PASSED ON FROM diseased PARENT to Children

Mitochondria Absent From Sperm "Head"

Maternal Inheritance of Mitochondrial DNA

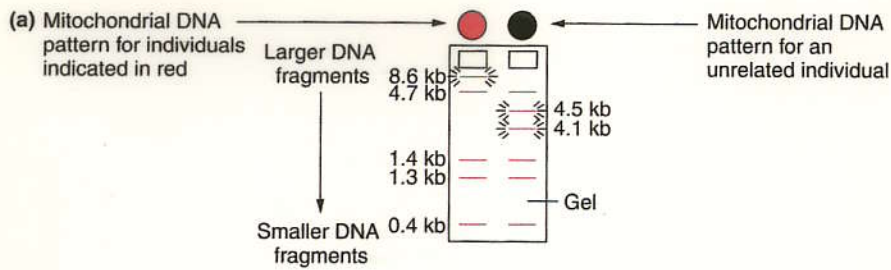
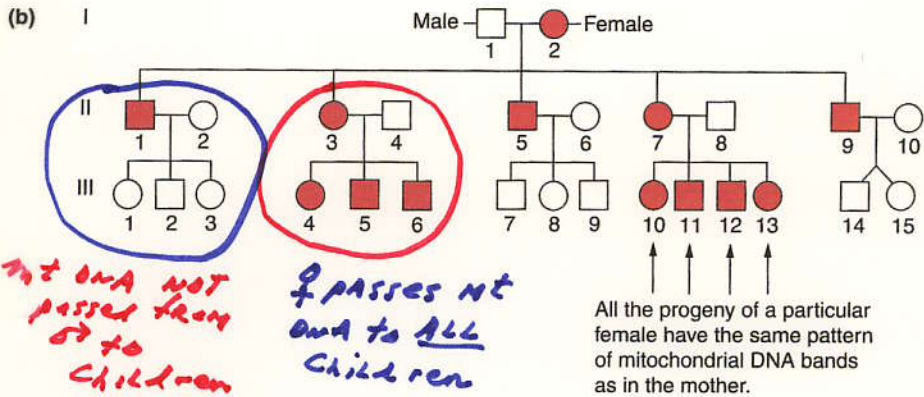
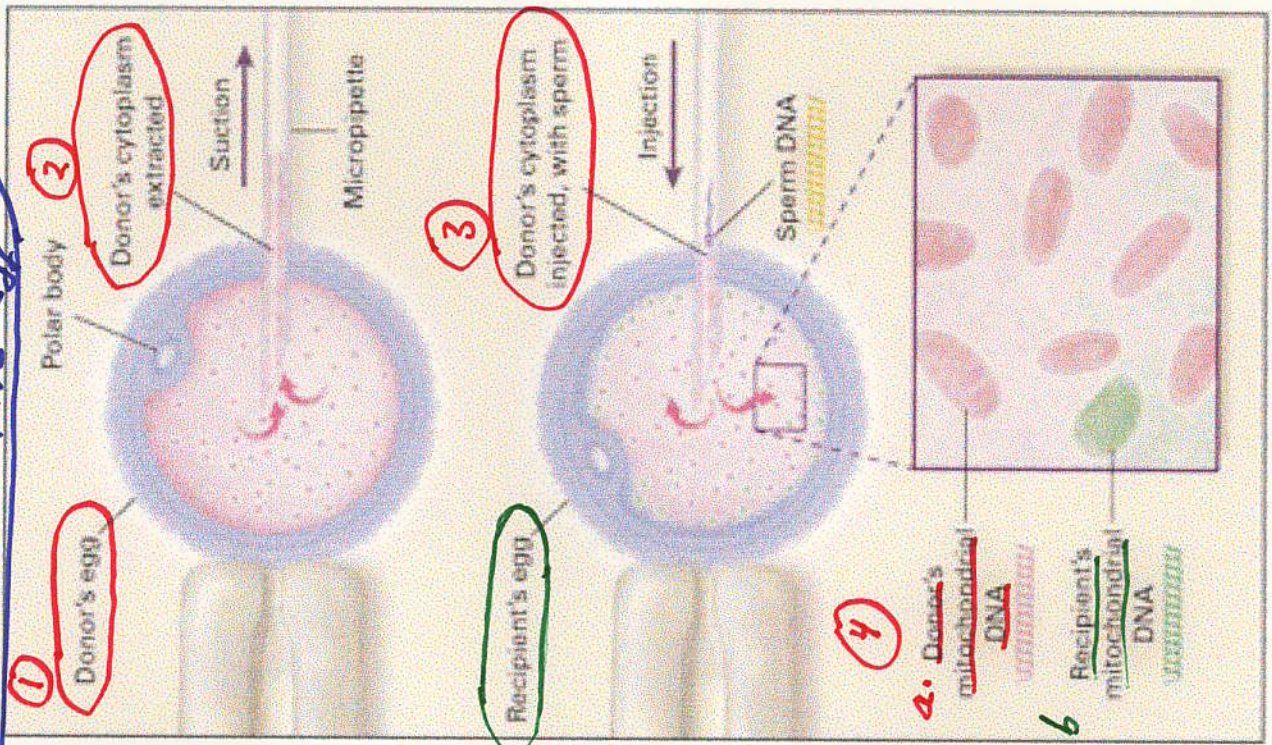


Figure 8.11 Maternal DNA Pattern of Inheritance Key genes for cell respiration and mitochondrial function are located in a small DNA ring in human mitochondria. Because mitochondria are contributed by the egg before fertilization, DNA can be traced through the maternal line with fingerprinting of the mitochondrial DNA.



**OOPASMIC TRANSFER TO INJECT
HEALTHY Mitochondrial Genomes
into Egg**



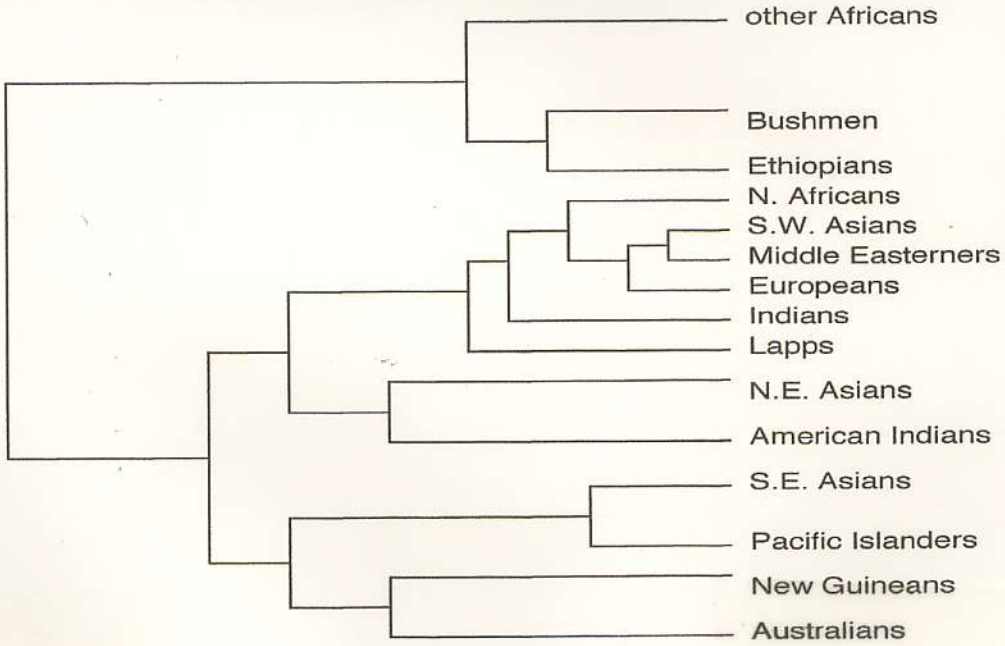
FORM OF
Gene Therapy
That's Inherited

MOSAIC!

Complements
M2. Gene
Defect
in Recipient's
Egg

USING MT DNA Polymorphisms to
construct trees of
Human origins

Evolutionary Tree



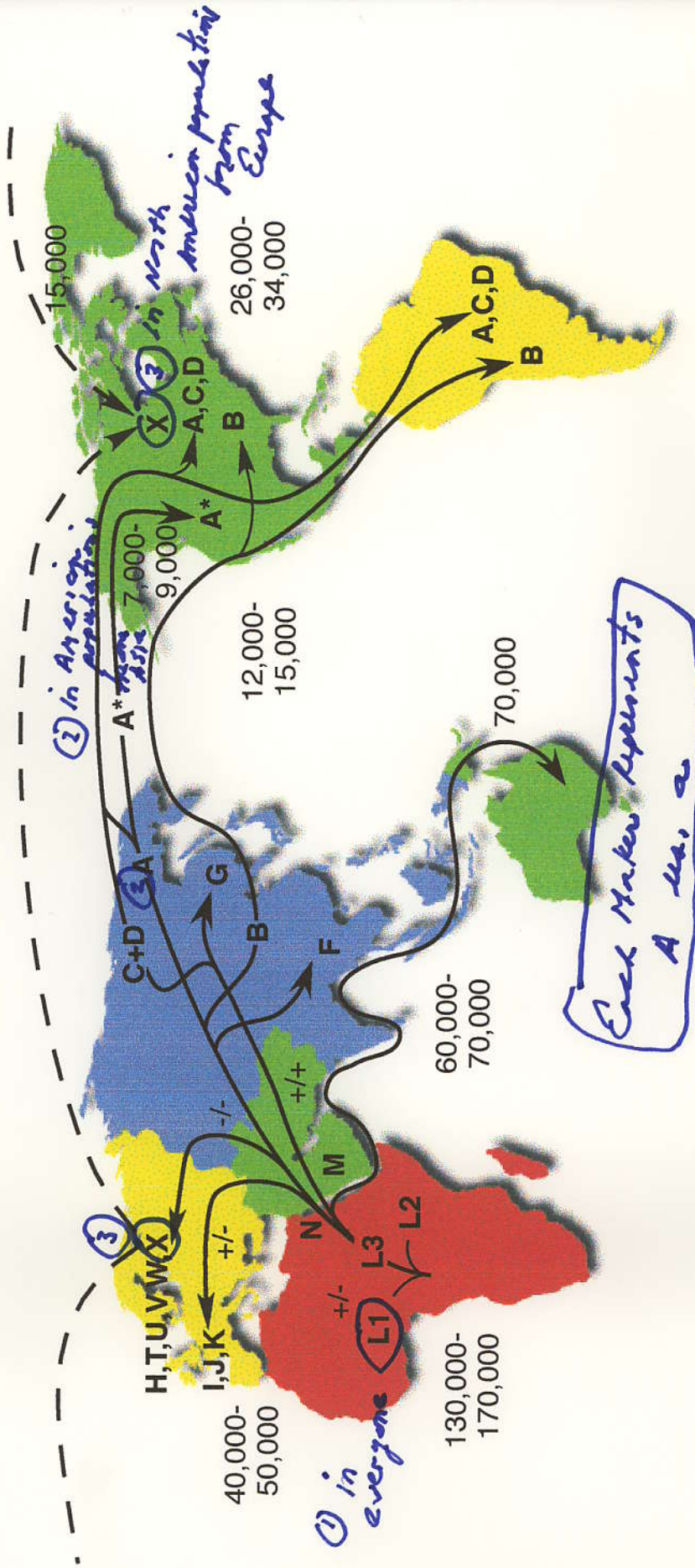
Passed on FROM "Eve" to us!

USING MT DNA RFLPs & SNPs

Human mtDNA Migrations

<http://www.mitomap.org/mitomap/WorldMigrations.pdf>

Copyright 2002 © Mitomap.org



+/-, +/+, or +/- = Dde I 10394 / Alu I 10397

* = Rsa I 16329

Mutation rate = 2.2 - 2.9 % / MYR

Time estimates are YBP

- 1** All populations will contain the oldest/ancient alleles!
- 2** New alleles/haplotypes arise in populations as they move/migrate & are fixed in the founder population

THE HUMAN GENOME SEQUENCE

articles

Initial sequencing and analysis of the human genome

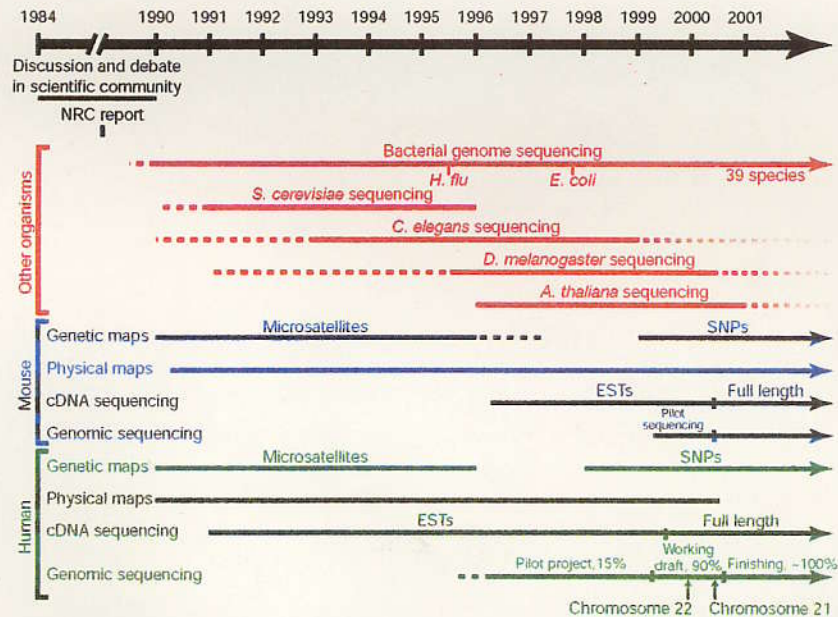
International Human Genome Sequencing Consortium*

* A partial list of authors appears on the opposite page. Affiliations are listed at the end of the paper.

The human genome holds an extraordinary trove of information about human development, physiology, medicine and evolution. Here we report the results of an international collaboration to produce and make freely available a draft sequence of the human genome. We also present an initial analysis of the data, describing some of the insights that can be gleaned from the sequence.

© 2001 Macmillan Magazines Ltd

NATURE | VOL 409 | 15 FEBRUARY 2001 | www.nature.com



Done!

Figure 1 Timeline of large-scale genomic analyses. Shown are selected components of work on several non-vertebrate model organisms (red), the mouse (blue) and the human (green) from 1990; earlier projects are described in the text. SNPs, single nucleotide polymorphisms; ESTs, expressed sequence tags.

13

THE HUMAN GENOME SEQUENCE IS THE RESULT OF AN INTERNATIONAL COLLABORATION

articles

Genome Sequencing Centres (Listed in order of total genomic sequence contributed, with a partial list of personnel. A full list of contributors at each centre is available as Supplementary Information.)

Whitehead Institute for Biomedical Research, Center for Genome Research: Eric S. Lander^{1*}, Lauren M. Linton¹, Bruce Birren^{1*}, Chad Nusbaum^{1*}, Michael C. Zody^{1*}, Jennifer Baldwin¹, Kerri Devon¹, Ken Dewar¹, Michael Doyle¹, William FitzHugh^{1*}, Roel Funke¹, Diane Gage¹, Katrina Harris¹, Andrew Heaford¹, John Howland¹, Lisa Kann¹, Jessica Lehoczyk¹, Rosie LeVine¹, Paul McEwan¹, Kevin McKernan¹, James Meldrim¹, Jill P. Mesirov^{1*}, Cher Miranda¹, William Morris¹, Jerome Naylor¹, Christina Raymond¹, Mark Rosetti¹, Ralph Santos¹, Andrew Sheridan¹, Carrie Sougnéz¹, Nicole Stange-Thomann¹, Nikola Stojanovic¹, Aravind Subramanian¹ & Dudley Wyman¹

The Sanger Centre: Jane Rogers², John Sulston^{2*}, Rachael Ainscough², Stephan Beck², David Bentley², John Burton², Christopher Clee², Nigel Carter², Alan Coulson², Rebecca Deadman², Panos Deloukas², Andrew Dunham², Ian Dunham², Richard Durbin^{2*}, Lisa French², Darren Grafham², Simon Gregory², Tim Hubbard^{2*}, Sean Humphray², Adrienne Hunt², Matthew Jones², Christine Lloyd², Amanda McMurray², Lucy Matthews², Simon Mercer², Sarah Milne², James C. Mullikin^{2*}, Andrew Mungall², Robert Plumb², Mark Ross², Ratna Showkeen² & Sarah Sims²

Washington University Genome Sequencing Center: Robert H. Waterston^{3*}, Richard K. Wilson³, LaDeana W. Hillier^{3*}, John D. McPherson³, Marco A. Marra³, Elaine R. Mardis³, Lucinda A. Fulton³, Asif T. Chhinwalla^{3*}, Kymberlie H. Pepin³, Warren R. Gish³, Stephanie L. Chissole³, Michael C. Wendt³, Kim D. Delehaunty³, Tracie L. Milner³, Andrew Delehaunty³, Jason B. Kramer³, Lisa L. Cook³, Robert S. Fulton³, Douglas L. Johnson³, Patrick J. Minx³ & Sandra W. Clifton³

US DOE Joint Genome Institute: Trevor Hawkins⁴, Elbert Branscomb⁴, Paul Predki⁴, Paul Richardson⁴, Sarah Wenning⁴, Tom Slezak⁴, Norman Doggett⁴, Jan-Fang Cheng⁴, Anne Olsen⁴, Susan Lucas⁴, Christopher Elkin⁴, Edward Uberbacher⁴ & Marvin Frazier⁴

Baylor College of Medicine Human Genome Sequencing Center: Richard A. Gibbs^{5*}, Donna M. Muzny⁵, Steven E. Scherer⁵, John B. Bouck^{5*}, Erica J. Sodergren⁵, Kim C. Worley^{5*}, Catherine M. Rives⁵, James H. Gorrell⁵, Michael L. Metzker⁵, Susan L. Naylor⁵, Raju S. Kucherlapati⁵, David L. Nelson⁵ & George M. Weinstock⁵

RIKEN Genomic Sciences Center: Yoshlyuki Sakaki⁶, Asao Fujiyama⁶, Masahira Hattori⁶, Tetsushi Yada⁶, Atsushi Toyoda⁶, Takehiko Itoh⁶, Chiharu Kawagoe⁶, Hidemi Watanabe⁶, Yasushi Totoki⁶ & Todd Taylor⁶

Genoscope and CNRS UMR-8030: Jean Weissenbach¹⁰, Roland Hellig¹⁰, William Saurin¹⁰, Francois Artiguenave¹⁰, Philippe Brottier¹⁰, Thomas Bruls¹⁰, Eric Pelletier¹⁰, Catherine Robert¹⁰ & Patrick Wincker¹⁰

GTC Sequencing Center: Douglas R. Smith¹¹, Lynn Doucette-Stamm¹¹, Marc Rubenfield¹¹, Keith Weinstock¹¹, Hong Mei Lee¹¹ & JoAnn Dubois¹¹

Department of Genome Analysis, Institute of Molecular

Biotechnology: André Rosenthal¹², Matthias Platzer¹², Gerald Nyakatura¹², Stefan Taudien¹² & Andreas Rump¹²

Beijing Genomics Institute/Human Genome Center: Huanming Yang¹³, Jun Yu¹³, Jian Wang¹³, Guyang Huang¹⁴ & Jun Gu¹⁵

Multimegabase Sequencing Center, The Institute for Systems Biology: Leroy Hood¹⁶, Lee Rowen¹⁶, Anup Madan¹⁶ & Shizen Qin¹⁶

Stanford Genome Technology Center: Ronald W. Davis¹⁷, Nancy A. Federspiel¹⁷, A. Pia Abola¹⁷ & Michael J. Proctor¹⁷

Stanford Human Genome Center: Richard M. Myers¹⁸, Jeremy Schmutz¹⁸, Mark Dickson¹⁸, Jane Grimwood¹⁸ & David R. Cox¹⁸

University of Washington Genome Center: Maynard V. Olson¹⁹, Rajinder Kaul¹⁹ & Christopher Raymond¹⁹

Department of Molecular Biology, Kelo University School of Medicine: Nobuyoshi Shimizu²⁰, Kazuhiko Kawasaki²⁰ & Shinsei Minoshima²⁰

University of Texas Southwestern Medical Center at Dallas: Glen A. Evans^{21†}, Maria Athanasiou²¹ & Roger Schultz²¹

University of Oklahoma's Advanced Center for Genome Technology: Bruce A. Roe²², Feng Chen²² & Huaqin Pan²²

Max Planck Institute for Molecular Genetics: Juliane Ramser²³, Hans Lehrach²³ & Richard Reinhardt²³

Cold Spring Harbor Laboratory, Lita Annenberg Hazen Genome Center: W. Richard McCombie²⁴, Melissa de la Bastide²⁴ & Neilay Dedhia²⁴

GBF—German Research Centre for Biotechnology: Helmut Blöcker²⁵, Klaus Hornischer²⁵ & Gabriele Nordtsiek²⁵

*** Genome Analysis Group (listed in alphabetical order, also includes individuals listed under other headings):** Richa Agarwala²⁶, L. Aravind²⁶, Jeffrey A. Bailey²⁷, Alex Bateman², Serafim Batzoglou¹, Ewan Birney²⁸, Peer Bork^{29,30}, Daniel G. Brown¹, Christopher B. Burge³¹, Lorenzo Cerutti²⁸, Hsiu-Chuan Chen²⁸, Deanna Church²⁸, Michele Clamp², Richard R. Copley³⁰, Tobias Doerks^{29,30}, Sean R. Eddy³², Evan E. Eichler²⁷, Terrence S. Furey²³, James Galagan¹, James G. R. Gilbert², Cyrus Harmon³⁴, Yoshitake Hayashizaki³⁵, David Haussler³⁶, Henning Hermjakob²⁸, Karsten Hokamp³⁷, Wonhee Jang²⁶, L. Steven Johnson³², Thomas A. Jones³², Simon Kasif³⁸, Arek Kasprzyk²⁸, Scot Kennedy³⁹, W. James Kent⁴⁰, Paul Kitts²⁶, Eugene V. Koonin²⁶, Ian Korf³, David Kulp³⁴, Doron Lancet⁴¹, Todd M. Lowe⁴², Aoife McLysaght³⁷, Tarjei Mikkelsen³⁵, John V. Moran⁴³, Nicola Mulder²⁸, Victor J. Pollar¹, Chris P. Ponting⁴⁴, Greg Schuler²⁶, Jörg Schultz³⁰, Guy Slater²⁸, Arian F. A. Smit⁴⁵, Ella Stupka²⁸, Joseph Szustakowski³⁸, Danielle Thierry-Mieg²⁶, Jean Thierry-Mieg²⁶, Lukas Wagner²⁶, John Wallis³, Raymond Wheeler²⁴, Alan Williams³⁴, Yuri I. Wolf²⁶, Kenneth H. Wolfe³⁷, Shlaw-Pyng Yang³ & Ru-Fang Yeh³¹

Scientific management: National Human Genome Research Institute, US National Institutes of Health: Francis Collins^{46*}, Mark S. Guyer⁴⁶, Jane Peterson⁴⁶, Adam Felsenfeld^{46*} & Kris A. Wetterstrand⁴⁶; Office of Science, US Department of Energy: Aristides Patrino⁴⁷; The Wellcome Trust: Michael J. Morgan⁴⁸

BUT IT WAS ALSO DONE INDEPENDENTLY
By A COMPANY - CELERA®

The Sequence of the Human Genome

J. Craig Venter,^{1*} Mark D. Adams,¹ Eugene W. Myers,¹ Peter W. Li,¹ Richard J. Mural,¹
Granger G. Sutton,¹ Hamilton O. Smith,¹ Mark Yandell,¹ Cheryl A. Evans,¹ Robert A. Holt,¹
Jeannine D. Gocayne,¹ Peter Amanatides,¹ Richard M. Ballew,¹ Daniel H. Huson,¹
Jennifer Russo Wortman,¹ Qing Zhang,¹ Chinnappa D. Kodira,¹ Xiangqun H. Zheng,¹ Lin Chen,¹
Marian Skupski,¹ Gangadharan Subramanian,¹ Paul D. Thomas,¹ Jinghui Zhang,¹
George L. Gabor Miklos,² Catherine Nelson,³ Samuel Broder,¹ Andrew G. Clark,⁴ Joe Nadeau,⁵
Victor A. McKusick,⁶ Norton Zinder,⁷ Arnold J. Levine,⁷ Richard J. Roberts,⁸ Mel Simon,⁹
Carolyn Slayman,¹⁰ Michael Hunkapiller,¹¹ Randall Bolanos,¹ Arthur Delcher,¹ Ian Dew,¹ Daniel Fasulo,¹
Michael Flanigan,¹ Lilliana Florea,¹ Aaron Halpern,¹ Sridhar Hannenhalli,¹ Saul Kravitz,¹ Samuel Levy,¹
Clark Mobarry,¹ Knut Reinert,¹ Karin Remington,¹ Jane Abu-Threideh,¹ Ellen Beasley,¹ Kendra Biddick,¹
Vivien Bonazzi,¹ Rhonda Brandon,¹ Michele Cargill,¹ Ishwar Chandramouliswaran,¹ Rosane Charlab,¹
Kabir Chaturvedi,¹ Zuoming Deng,¹ Valentina Di Francesco,¹ Patrick Dunn,¹ Karen Eilbeck,¹
Carlos Evangelista,¹ Andrei E. Gabrielian,¹ Weiniu Gan,¹ Wangmao Ge,¹ Fangcheng Gong,¹ Zhiping Gu,¹
Ping Guan,¹ Thomas J. Heiman,¹ Maureen E. Higgins,¹ Rui-Ru Ji,¹ Zhaoxi Ke,¹ Karen A. Ketchum,¹
Zhongwu Lai,¹ Yiding Lei,¹ Zhenya Li,¹ Jiayin Li,¹ Yong Liang,¹ Xiaoying Lin,¹ Fu Lu,¹
Gennady V. Merkulov,¹ Natalia Milshina,¹ Helen M. Moore,¹ Ashwinikumar K Naik,¹
Vaibhav A. Narayan,¹ Beena Neelam,¹ Deborah Nusskern,¹ Douglas B. Rusch,¹ Steven Salzberg,¹²
Wei Shao,¹ Bixiong Shue,¹ Jingtao Sun,¹ Zhen Yuan Wang,¹ Aihui Wang,¹ Xin Wang,¹ Jian Wang,¹
Ming-Hui Wei,¹ Ron Wides,¹³ Chunlin Xiao,¹ Chunhua Yan,¹ Alison Yao,¹ Jane Ye,¹ Ming Zhan,¹
Weiqing Zhang,¹ Hongyu Zhang,¹ Qi Zhao,¹ Liansheng Zheng,¹ Fei Zhong,¹ Wenyan Zhong,¹
Shiaoping C. Zhu,¹ Shaying Zhao,¹² Dennis Gilbert,¹ Suzanna Baumhueter,¹ Gene Spier,¹
Christine Carter,¹ Anibal Cravchik,¹ Trevor Woodage,¹ Feroze Ali,¹ Huijin An,¹ Aderonke Awe,¹
Danita Baldwin,¹ Holly Baden,¹ Mary Barnstead,¹ Ian Barrow,¹ Karen Beeson,¹ Dana Busam,¹
Amy Carver,¹ Angela Center,¹ Ming Lai Cheng,¹ Liz Curry,¹ Steve Danaher,¹ Lionel Davenport,¹
Raymond Desilets,¹ Susanne Dietz,¹ Kristina Dodson,¹ Lisa Doup,¹ Steven Ferreira,¹ Neha Garg,¹
Andres Gluecksmann,¹ Brit Hart,¹ Jason Haynes,¹ Charles Haynes,¹ Cheryl Heiner,¹ Suzanne Hladun,¹
Damon Hostin,¹ Jarrett Houck,¹ Timothy Howland,¹ Chinyere Ibegwam,¹ Jeffery Johnson,¹
Francis Kalush,¹ Lesley Kline,¹ Shashi Koduru,¹ Amy Love,¹ Felecia Mann,¹ David May,¹
Steven McCawley,¹ Tina McIntosh,¹ Ivy McMullen,¹ Mee Moy,¹ Linda Moy,¹ Brian Murphy,¹
Keith Nelson,¹ Cynthia Pfannkoch,¹ Eric Pratts,¹ Vinita Puri,¹ Hina Qureshi,¹ Matthew Reardon,¹
Robert Rodriguez,¹ Yu-Hui Rogers,¹ Deanna Romblad,¹ Bob Ruhfel,¹ Richard Scott,¹ Cynthia Sitter,¹
Michelle Smallwood,¹ Erin Stewart,¹ Renee Strong,¹ Ellen Suh,¹ Reginald Thomas,¹ Ni Ni Tint,¹
Sukye Tse,¹ Claire Vech,¹ Gary Wang,¹ Jeremy Wetter,¹ Sherita Williams,¹ Monica Williams,¹
Sandra Windsor,¹ Emily Winn-Deen,¹ Keriellen Wolfe,¹ Jayshree Zaveri,¹ Karena Zaveri,¹
Josep F. Abril,¹⁴ Roderic Guigó,¹⁴ Michael J. Campbell,¹ Kimmen V. Sjolander,¹ Brian Karlak,¹
Anish Kejariwal,¹ Huaiyu Mi,¹ Betty Lazareva,¹ Thomas Hatton,¹ Apurva Narechania,¹ Karen Diemer,¹
Anushya Muruganujan,¹ Nan Guo,¹ Shinji Sato,¹ Vineet Bafna,¹ Sorin Istrail,¹ Ross Lippert,¹
Russell Schwartz,¹ Brian Walenz,¹ Shibu Yooseph,¹ David Allen,¹ Anand Basu,¹ James Baxendale,¹
Louis Blick,¹ Marcelo Caminha,¹ John Carnes-Stine,¹ Parris Caulk,¹ Yen-Hui Chiang,¹ My Coyne,¹
Carl Dahlke,¹ Anne Deslattes Mays,¹ Maria Dombroski,¹ Michael Donnelly,¹ Dale Ely,¹ Shiva Esparham,¹
Carl Fosler,¹ Harold Gire,¹ Stephen Glanowski,¹ Kenneth Glasser,¹ Anna Glodek,¹ Mark Gorokhov,¹
Ken Graham,¹ Barry Gropman,¹ Michael Harris,¹ Jeremy Heil,¹ Scott Henderson,¹ Jeffrey Hoover,¹
Donald Jennings,¹ Catherine Jordan,¹ James Jordan,¹ John Kasha,¹ Leonid Kagan,¹ Cheryl Kraft,¹
Alexander Levitsky,¹ Mark Lewis,¹ Xiangjun Liu,¹ John Lopez,¹ Daniel Ma,¹ William Majoros,¹
Joe McDaniel,¹ Sean Murphy,¹ Matthew Newman,¹ Trung Nguyen,¹ Ngoc Nguyen,¹ Marc Nodell,¹
Sue Pan,¹ Jim Peck,¹ Marshall Peterson,¹ William Rowe,¹ Robert Sanders,¹ John Scott,¹
Michael Simpson,¹ Thomas Smith,¹ Arlan Sprague,¹ Timothy Stockwell,¹ Russell Turner,¹ Eli Venter,¹
Mei Wang,¹ Meiyuan Wen,¹ David Wu,¹ Mitchell Wu,¹ Ashley Xia,¹ Ali Zandieh,¹ Xiaohong Zhu¹

AND COMPLETED IN ONLY NINE MONTHS!

IT WAS A RACE!

Celera Effort

5 people's DNA
combined & sequenced

3 ♀ 2 ♂

- 2 Caucasian
- 1 African American
- 1 Chinese
- 1 Hispanic

↓ SNPs
= allele variability

2005
5.8 x 10⁶ SNPs
or
1/2600 bp
on average between
2 people

A 2.91-billion base pair (bp) consensus sequence of the euchromatic portion of the human genome was generated by the whole-genome shotgun sequencing method. The 14.8-billion bp DNA sequence was generated over 9 months from 27,271,853 high-quality sequence reads (5.11-fold coverage of the genome) from both ends of plasmid clones made from the DNA of five individuals. Two assembly strategies—a whole-genome assembly and a regional chromosome assembly—were used, each combining sequence data from Celera and the publicly funded genome effort. The public data were shredded into 550-bp segments to create a 2.9-fold coverage of those genome regions that had been sequenced, without including biases inherent in the cloning and assembly procedure used by the publicly funded group. This brought the effective coverage in the assemblies to eightfold, reducing the number and size of gaps in the final assembly over what would be obtained with 5.11-fold coverage. The two assembly strategies yielded very similar results that largely agree with independent mapping data. The assemblies effectively cover the euchromatic regions of the human chromosomes. More than 90% of the genome is in scaffold assemblies of 100,000 bp or more, and 25% of the genome is in scaffolds of 10 million bp or larger. Analysis of the genome sequence revealed 26,588 protein-encoding transcripts for which there was strong corroborating evidence and an additional ~12,000 computationally derived genes with mouse matches or other weak supporting evidence. Although gene-dense clusters are obvious, almost half the genes are dispersed in low G+C sequence separated by large tracts of apparently noncoding sequence. Only 1.1% of the genome is spanned by exons, whereas 24% is in introns, with 75% of the genome being intergenic DNA. Duplications of segmental blocks, ranging in size up to chromosomal lengths, are abundant throughout the genome and reveal a complex evolutionary history. Comparative genomic analysis indicates vertebrate expansions of genes associated with neuronal function, with tissue-specific developmental regulation, and with the hemostasis and immune systems. DNA sequence comparisons between the consensus sequence and publicly funded genome data provided locations of 2.1 million single-nucleotide polymorphisms (SNPs). A random pair of human haploid genomes differed at a rate of 1 bp per 1250 on average, but there was marked heterogeneity in the level of polymorphism across the genome. Less than 1% of all SNPs resulted in variation in proteins, but the task of determining which SNPs have functional consequences remains an open challenge.

1% Exons
24% Introns
75% Intergenic
100% DNA

1 SNP / 1250 bp on average between two people

Now one SNP / 1000 bp - HAPMAP PROJECT

BUT The International Public Sequence is More Complete than the Private one That contains Gaps

READ -

"The GENOME WAR"

by James Shreeve

ISBN 0375406298

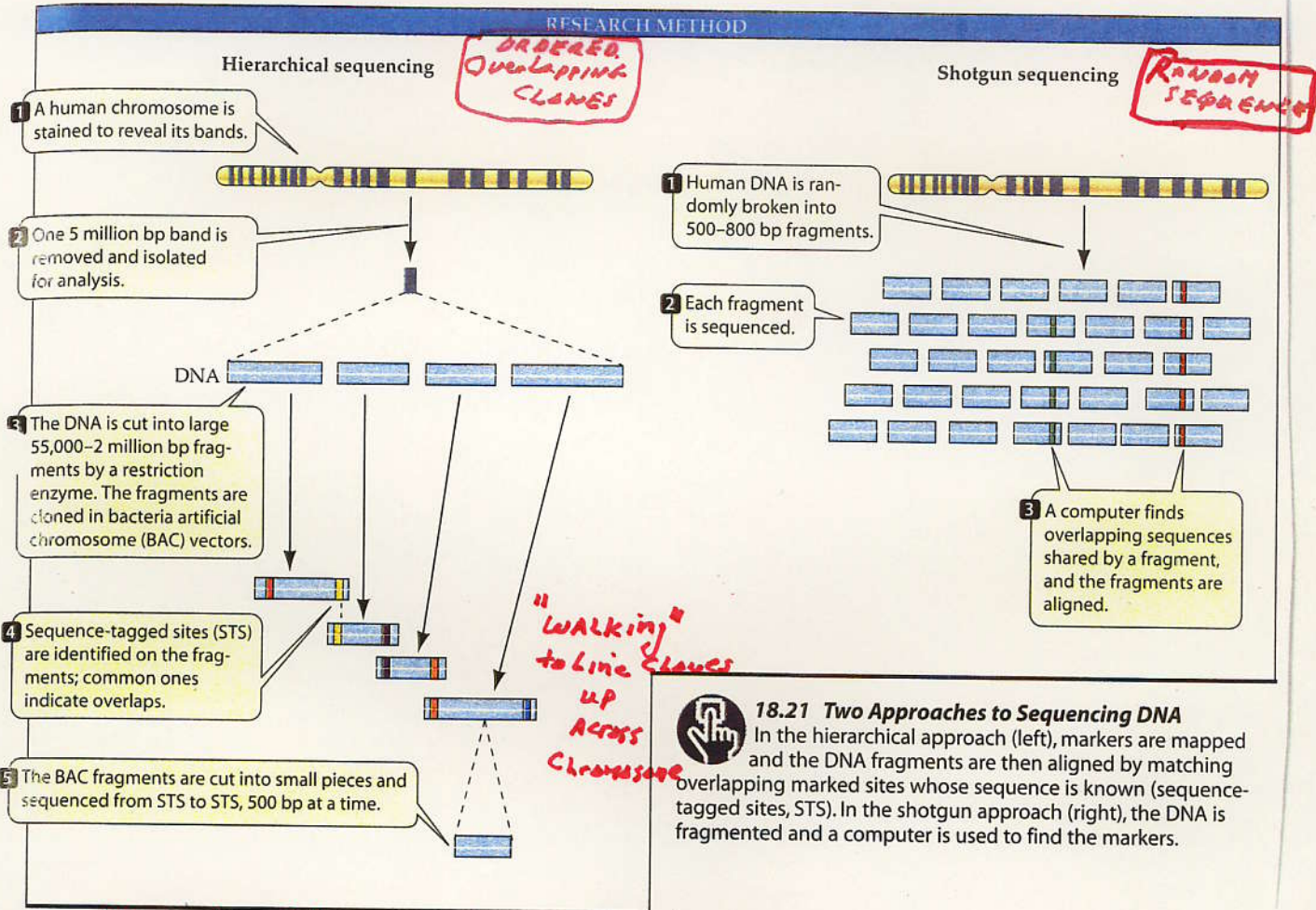
JANUARY, 2004

Private vs. Public Genome Projects!

HOW WAS THE HUMAN GENOME SEQUENCED?

TOP DOWN
Physical Map → Sequence

BOTTOM UP
Sequence → Assemble



18.21 Two Approaches to Sequencing DNA
In the hierarchical approach (left), markers are mapped and the DNA fragments are then aligned by matching overlapping marked sites whose sequence is known (sequence-tagged sites, STS). In the shotgun approach (right), the DNA is fragmented and a computer is used to find the markers.

23 contigs (8) 24 contigs (5) - Y chromosome

PUBLIC EFFORT

CHROMOSOME WALKS
↓
SEQUENCE

BEST
MOST COMPLETE
ENTIRE
SLOW

PRIVATE EFFORT

SHOTGUN

FAST
MANY GAPS
SKELETON

NEEDS PUBLIC DATA
TO ASSEMBLE → CHANGES

18

Approximately Half of the Genes
in the Human Genome have
Unknown Functions

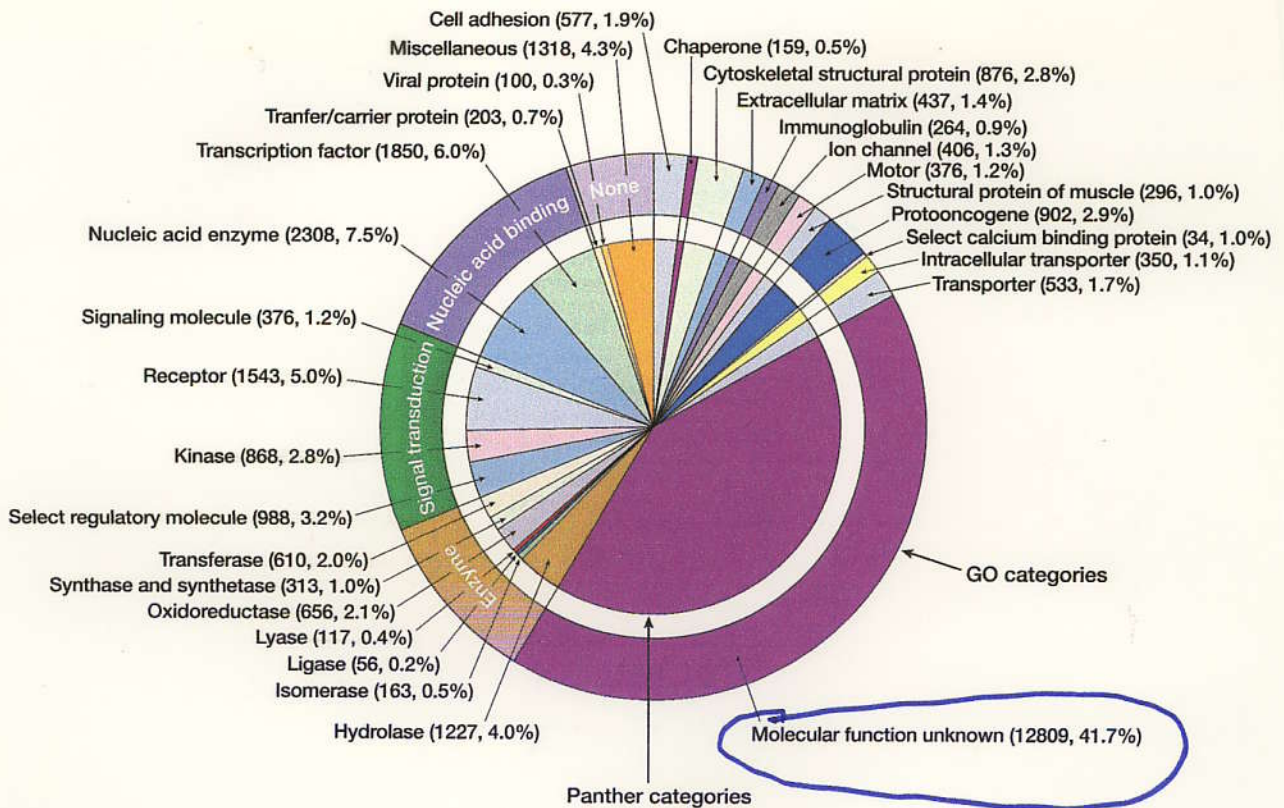


Figure 12.20: A preliminary functional classification of human polypeptide-encoding genes.

Known or predicted functions for 26 383 human polypeptide-encoding genes. Classification is according to the GO molecular function categories as shown in the outer circle (Gene Ontology classification – see Section 8.3.6) or to Celera's Panther molecular function categories (inner circle). Reproduced from Venter *et al.* (2001) *Science* **291**, 1304–1351, with permission from the American Association for the Advancement of Science.

LOTS OF WORK
yet to do!

Use Mouse as Model to Identify
Unknown Genes Functions!!

VERTABRATE AND MAMMALION RELATIONSHIPS

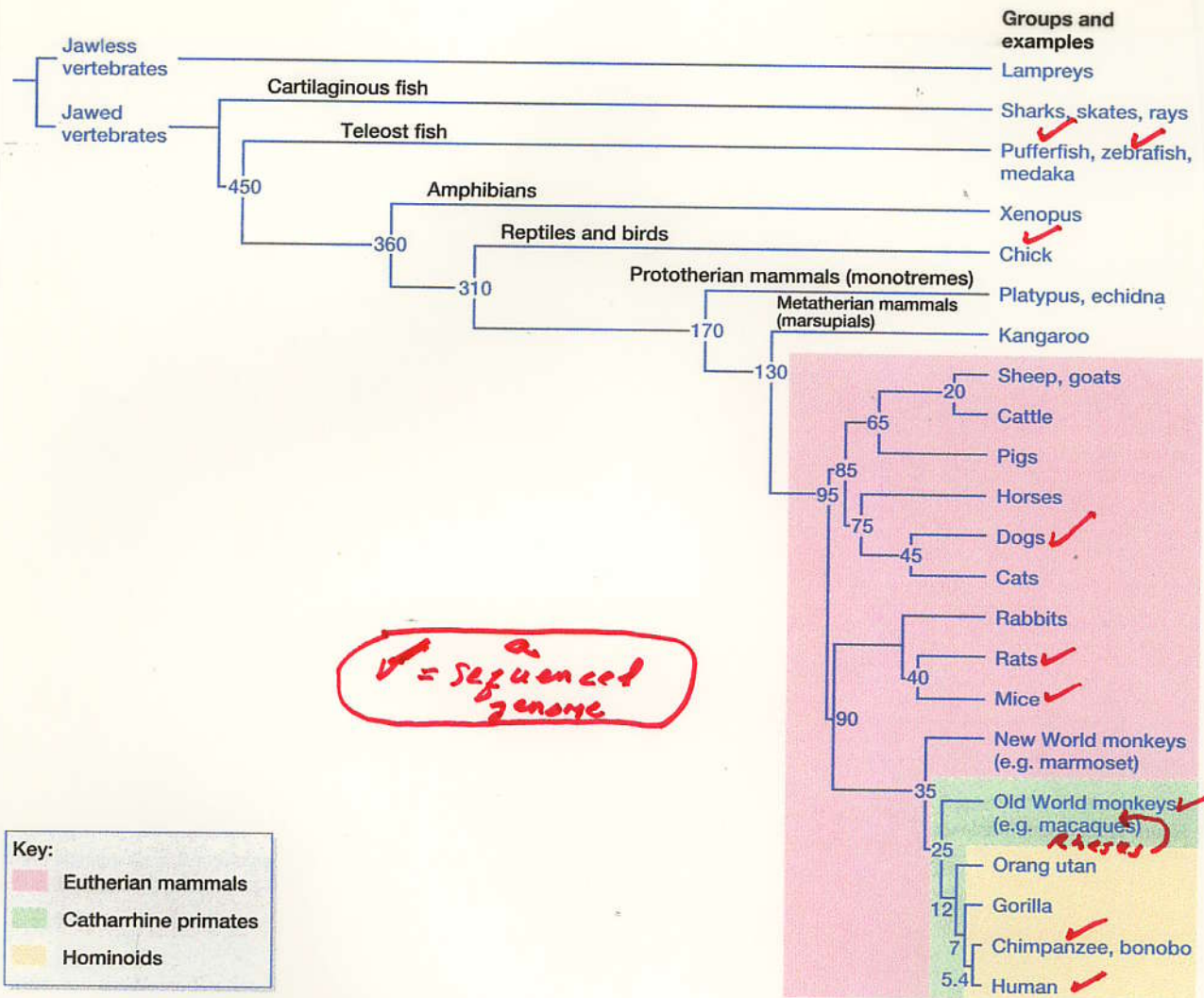


Figure 12.24: A simplified vertebrate phylogeny.
 Numbers at nodes show estimated divergence times in millions of years.

COMPARATIVE GENOMICS

HUMAN & CHIMP GENOMES ARE
VIRTUALLY IDENTICAL

Chimp Genome Draft Online

The relationship between humans and chimps just got a little easier to understand. This week, the consortium that has been unraveling the DNA sequence of our closest cousin for the past year put its results into the public domain. Robert Waterston of the University of Washington, Seattle, and his colleagues at the Broad Institute in Cambridge, Massachusetts (including the former genome center of the Whitehead Institute for Biomedical Research), and at the Washington University Genome Sequencing Center in St. Louis, Missouri, have determined the order of many of the 3.1 billion bases of a single male chimp's genome. Until now, only pieces of deciphered DNA were available, and their placement along the two dozen chromosomes was uncertain.

On average, each base was sequenced only four times. That's far

short of the current 10-fold coverage of the human genome, but it's enough to put together a rough draft with many bases in the right order, which the researchers deposited online in GenBank. The consortium matched up the human and chimp genomes base by base as much as possible, an alignment that will make it easier for researchers to find elusive genes and regulatory regions. The matchup will also highlight specific differences between the two genomes, perhaps further hinting at what set us apart. The sequence is more than researchers could have hoped for a few years ago, says Ajit Varki of the University of California, San Diego, and it "will be most useful" for geneticists trying to find genes responsible for inherited diseases.

Work on the chimp sequence continues, but in the meantime, the consortium is taking the next few months to analyze the data it has. It expects to publish results early next year.

-E.P.

CREDITS: (LEFT TO RIGHT) PAT POWERS AND CHERRY SCHAFER/GETTY IMAGES; E

10,000,000 YEAR DIVERGENCE

BUT differences should indicate why a
man is a man & a chimp a chimp

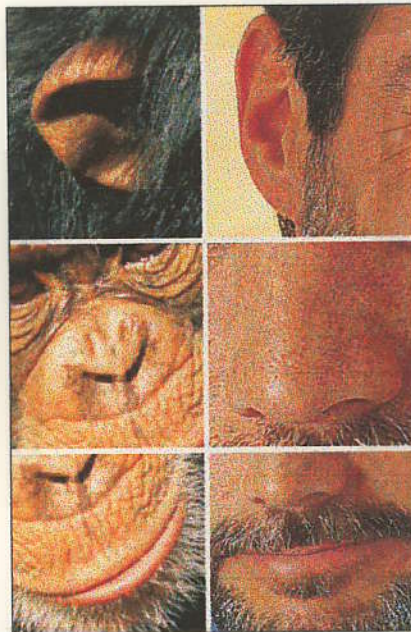
Key to understanding unique
Human Features at Molecular
Level!!

Compare all Mammalian Genomes!

NEWS OF THE WEEK

EVOLUTION

Genome Comparisons Hold Clues to Human Evolution



Hear no evil. Changes in genes for hearing, olfaction, and speech helped prompt human evolution.

SINGLE BASE PAIR CHANGES (SNPs)
ARE FREQUENT IN GENOME

TABLE 9.1 Five Classes of DNA Polymorphism

Class	Cause	Rate of Mutation per Locus per Gamete	Frequency in Genome	Number per Human Genome (on average)
Single base	Mutagens or replication errors	10^{-8} - 10^{-9}	1/700 bp	3 million
Microsatellite	Slippage during replication	10^{-3}	1/30,000 bp	100,000
Minisatellite	Unequal crossovers	10^{-3}	Unknown; discovered by chance	Fewer than 100 families known, yielding 1000 copies in all
Deletions	Mutagens; unequal crossovers	Extremely rare	Very low	0 - a few
Duplications	Mutagens; unequal crossovers	Extremely rare	Very low	0 - a few
Other insertions (excluding those resulting from micro- or minisatellite recombination)	Transposable elements	Extremely rare	Very low	0 - a few
Complex haplotype (any locus of 5 kb or more)	Any of the above	Combination of the above	Not applicable	Not applicable

DETECTED USING RESTRICTION ENZYMES (THE OLD FASHIONED WAY) OR BY DIRECT DNA SEQUENCING OF 2 INDIVIDUALS' GENES/GENOMES

SNPs or Single Nucleotide Polymorphisms occur randomly in the human genome

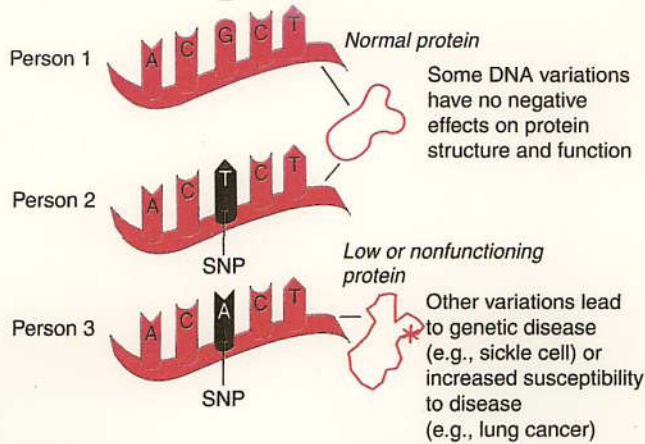


Figure 1.10 Single Nucleotide Polymorphisms A small piece of a gene sequence for three different individuals is represented. For simplicity, only one strand of a DNA molecule is shown. Notice how person 2 has a SNP in this gene which has no effect on protein structure and function. Person 3 however, has a different SNP in the same gene. This subtle genetic change may affect how this person responds to a medical drug or influence the likelihood that person 3 will develop a genetic disease.

Detected directly by sequencing

used to individualize genomes
*
associate with disease genes (i.e., disease gene markers)

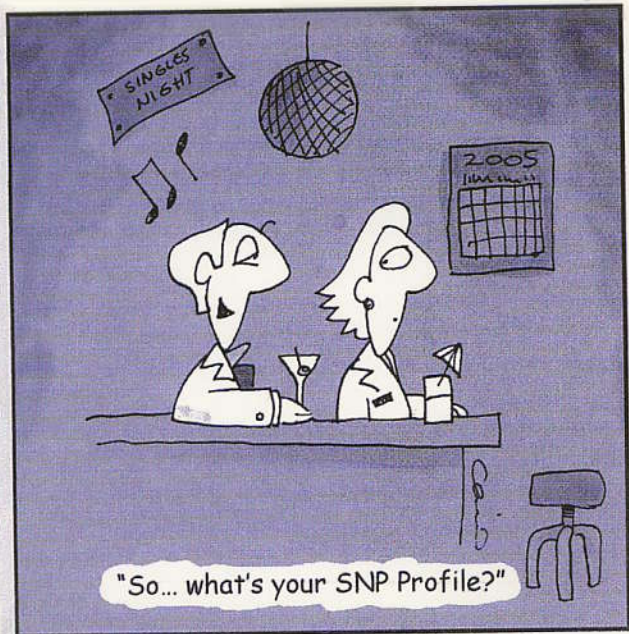


Figure 1.11 Secrets of the Human Genome In the future, we will have unprecedented knowledge of our genetic make-up including SNPs and other markers of genetic diseases. Can you think of possible ethical, legal, and social implications of such information?

LOCATION of SNPs Relative to Genes

TABLE 18-11 dbSNP Statistics (NCBI Genome Build 30, November 2002)

See http://www.ncbi.nlm.nih.gov/SNP/snp_summary.cgi: SNP count is the number of distinct RefSNPs having the noted functional relationship to at least one mRNA in the current assembly. Gene count is the number of distinct locus_id(s) having at least one variation of the noted functional class. (Genes with multiple variations may be counted in multiple classes.)

Functional Classification	SNP Count	Gene Count
Locus region	291,459	26,482
Allele synonymous to contig nucleotide	12,322	7,147
Allele nonsynonymous to contig nucleotide	16,251	8,496
Untranslated region	131,987	13,208
Intron ←	904,573	22,113
Splice site	277	268
Allele is same as contig nucleotide	28,491	11,621
Coding: synonymy unknown	13,501	3,584

TABLE 18-10 SNP Resources

Resource	Comment	URL
dbSNP	At NCBI	▶ www.ncbi.nlm.nih.gov/SNP
Human SNP database	At the Whitehead Institute	▶ http://www-genome.wi.mit.edu/snp/human/
The SNP Consortium (TSC)	A collaboration of industrial and academic laboratories	▶ http://snp.cshl.org

MOST SNPs ARE NOT IN CODING SEQUENCES & HAVE NO PHENOTYPIC EFFECTS!

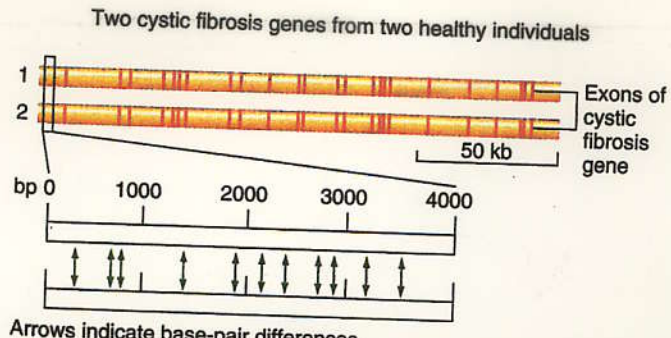


Figure 9.2 Base-pair differences between DNA cloned from the cystic fibrosis locus of two healthy individuals. These base-pair differences have no phenotypic effect; apparently they neither encode nor regulate expressed regions of the gene.

NOTE Differences Between Two Healthy Cystic Fibrosis Genes!

$\sim 6 \times 10^6$ SNPs (2006) !!
 $3.3 \times 10^9 / 6 \times 10^6 = 500-600 \text{ bp/SNP on average}$

EACH OF US DIFFERS BY $6 \times 10^6 \text{ bp}$ or $\sim 0.1\%$ of genome !!

MOST in NON-CODING/INTRON + Intergenic Regions & Useful As Markers For Forensics, Disease Genes, Populations

USING SNPs AS GENE MARKERS

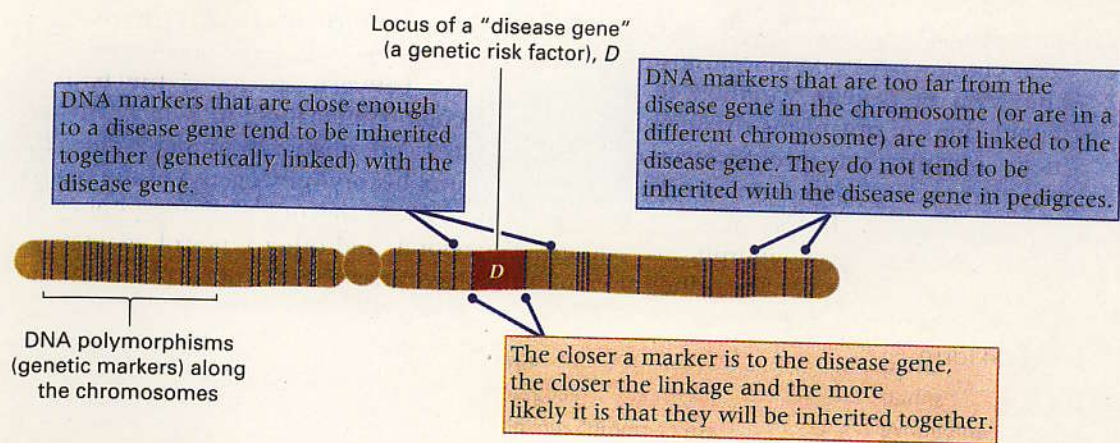


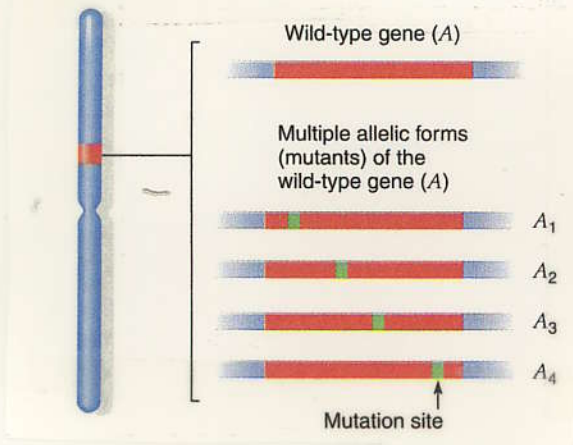
Figure 2.29 Concepts in genetic localization of genetic risk factors for disease. Polymorphic DNA markers (indicated by the vertical lines) that are close to a genetic risk factor (*D*) in the chromosome tend to be inherited together with the disease itself. The genomic location of the risk factor is determined by examining the known genomic locations of the DNA polymorphisms that are linked with it.

MAP GENES FOR DISEASE SUSCEPTIBILITY -
 GENES ENCODING COMPLEX (MULTI-GENIC)
 TRAITS (e.g., Heart Disease, Depression, Drug
 Sensitivity, Obesity)

SNPs = INDIVIDUAL GENE PROFILE
 ↳ Individual Medicine!

SNPs Generate Multiple Alleles
in a Population

Figure 4.1
Allelic forms of a gene.



Recall - 2 alleles/individual
Many in POPULATIONS!

USES OF SNPs

- ① Gene Identity / Allele Marker
- ② Disease Gene Identity / Pedigrees / Testing
- ③ Group Identity / Population History / Human origins
- ④ Individual Identity
 - ↳ Pharmacogenomics
 - ↳ Disease Prevention / Preventative Medicine
 - ↳ DNA Fingerprinting
- ⑤ Group Susceptibility to Disease / Drugs
 - ↳ Identify Genes Involved
 - ↳ "Resistance" Genes

Using SNPs to Personalize Medicine Pharmacogenomics

Individuals respond differently to the anti-leukemia drug 6-mercaptopurine.

Most people metabolize the drug quickly. Doses need to be high enough to treat leukemia and prevent relapses.



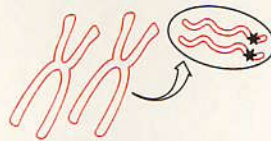
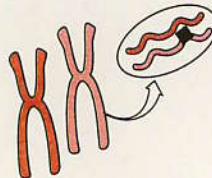
Others metabolize the drug slowly and need lower doses to avoid toxic side effects of the drug.



A small portion of people metabolize the drug so poorly that its effects can be fatal.



The diversity in responses is due to variations (mutations, ■ or *) in the gene for an enzyme called TPMT, or thiopurine methyltransferase.



After a simple blood test, individuals can be given doses of medication that are tailored to their genetic profile.



Normal dose

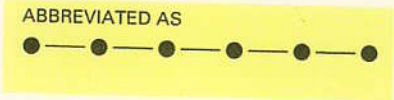
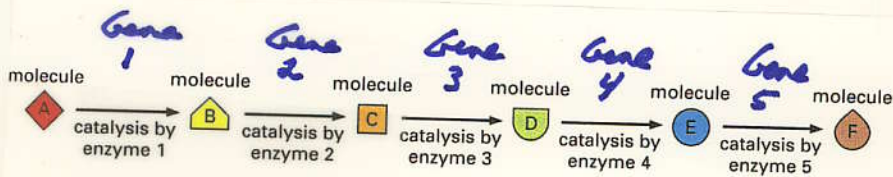


Dose for an extra slow metabolizer (TPMT deficient)



Figure 11.7 Pharmacogenomics Different individuals with the same disease often respond differently to a drug treatment because of subtle differences in gene expression. The dose that works for one person may be toxic for another—this is a basic problem of conventional medicine. Pharmacogenomics holds the promise of customizing medical treatment by determining the appropriate dosage for each individual based on the genes that person expresses.

Using SNPs in Pharmacogenetics



Different alleles can give rise to enzymes that differ slightly in activity.

Figure 2-34 How a set of enzyme-catalyzed reactions generates a metabolic pathway. Each enzyme catalyzes a particular chemical reaction, leaving the enzyme unchanged. In this example, a set of enzymes acting in series converts molecule A to molecule F, forming a metabolic pathway.

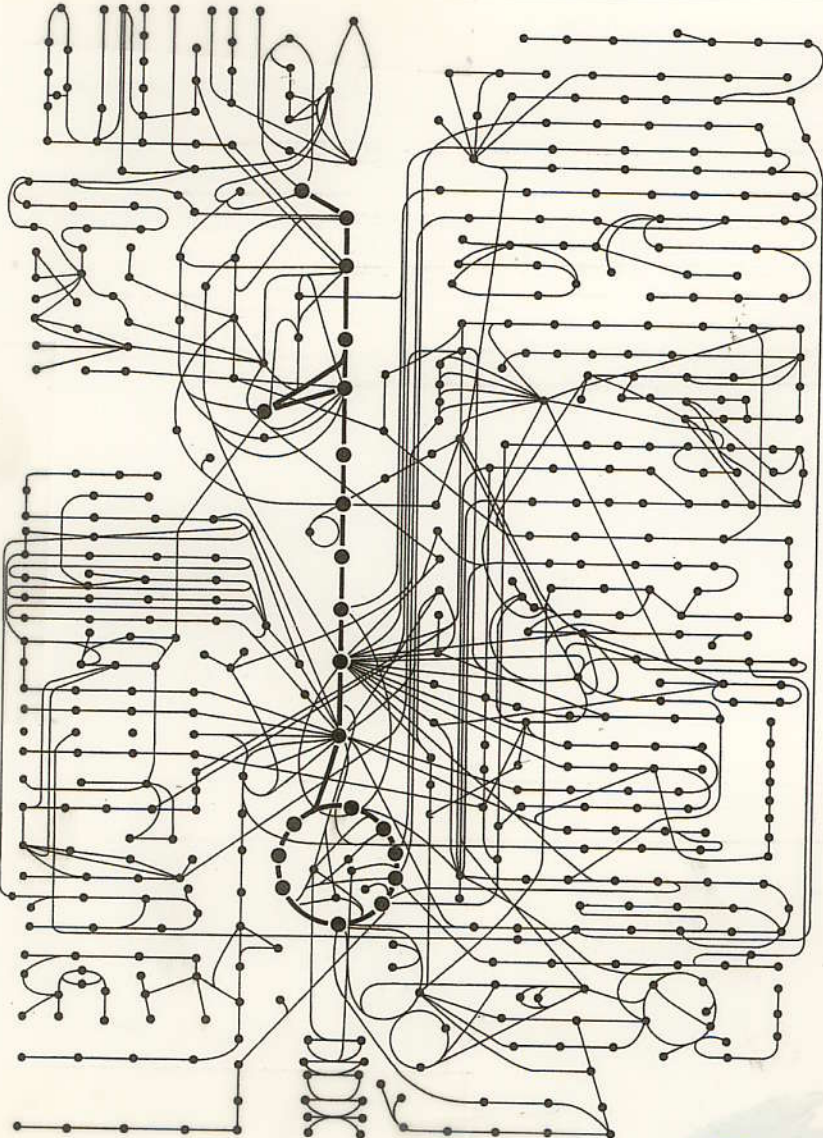


Figure 2-35 Some of the metabolic pathways and their interconnections in a typical cell. About 500 common metabolic reactions are shown diagrammatically, with each molecule in a metabolic pathway represented by a filled circle, as in the yellow box in Figure 2-34.

Goal: to find genes that encode enzymes that efficiently metabolize drugs
∴ no side effects

MOLECULAR SNP DETECTION

① RFLPs → BLOTS
→ PCR

② ASOs → PCR + Hybridization / "Blot"

③ Comparative Sequencing / Genes / Genome

④ Chips / Whole Genome / Individuals
+ Groups

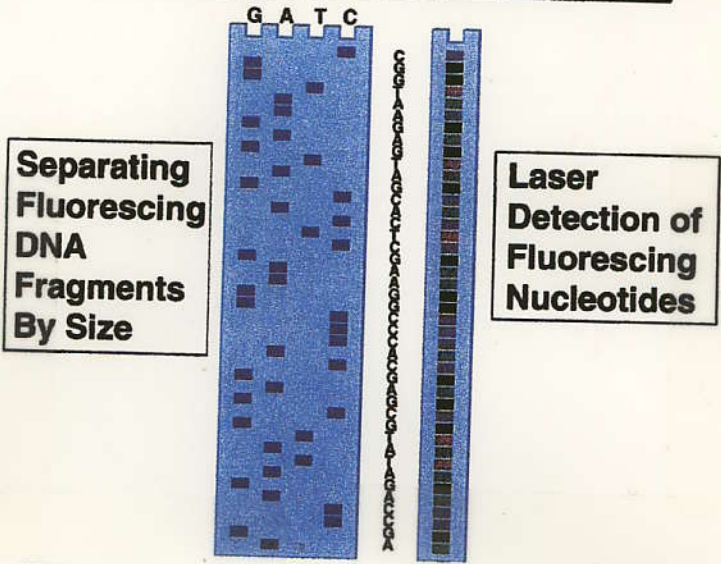
after SNPs identified by sequencing
+ synthesized as ASOs on chip!

DIRECT DNA SEQUENCING
+
CHIPS TO DETECT
SNPS

Gene & Whole Genome Approaches

Detecting SNPs By Sequencing

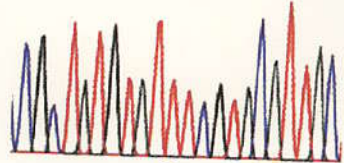
Genome Sequencing Using Computers and Robotics



Individual 1

27-D

CGC TATG TAT TCGTACATTAC
90 100

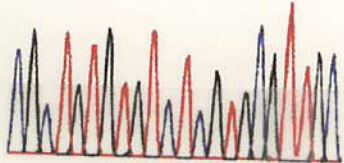


16093 T

Individual 2

28-D

CGC TATG TAT TCGTACATTAC
90 100



16093 C

Using Chips to Detect SNPs

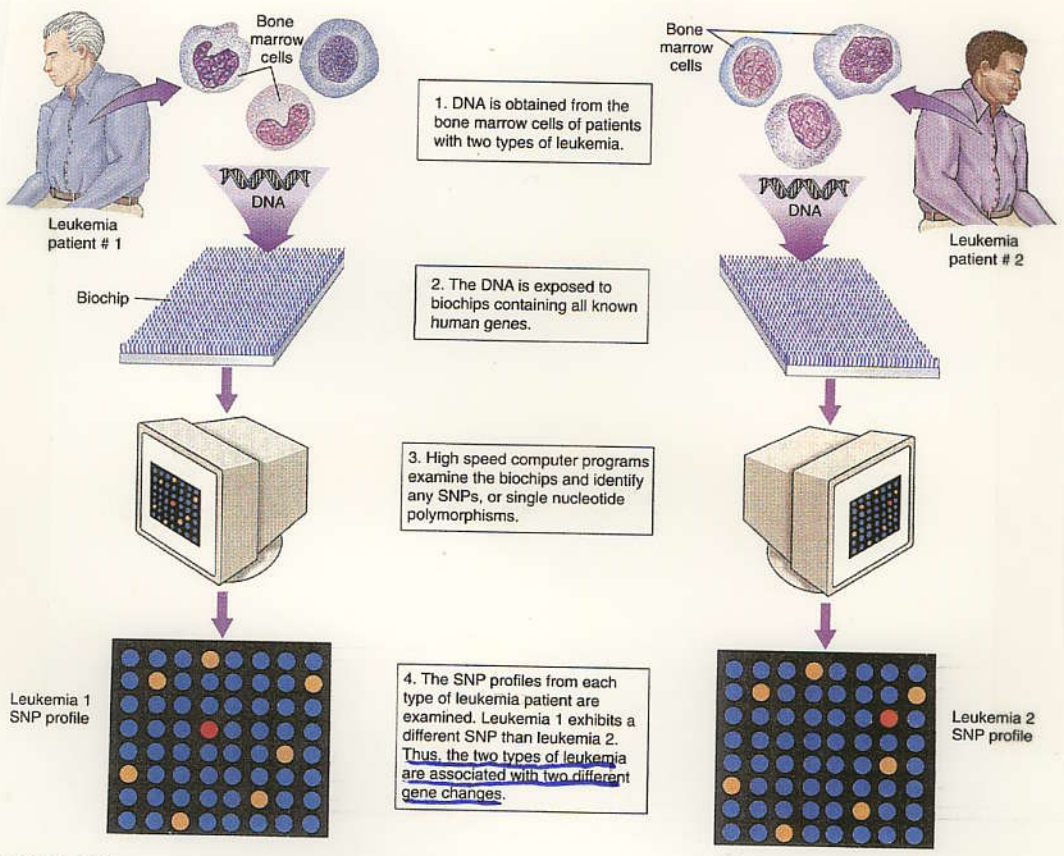


FIGURE 19.16 Biochips can help in identifying precise forms of cancer.

HAPLOTYPES OR
SNPs on a Chromosome
Inherited as a Unit

There are Millions of SNPs
That Differ Among
Individuals

. . . . But a small few reflect
our ancestry & "travel" in
groups on chromosomes —
are linked & may show
specific gene linkages!

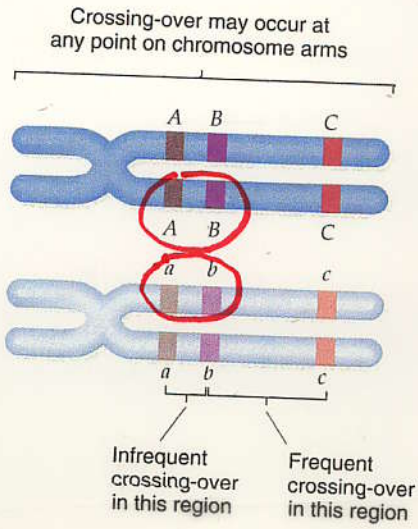
Haplotypes!

CLOSELY-LINKED SNPs ARE INHERITED AS A UNIT

Figure 5.3

The relationship between recombination and map distance.

The farther apart two genes are, the greater the number of possible sites for recombination. Thus, the probability of recombination occurring between genes A and B is much less than that between genes B and C. The percentage of recombinants can provide information about the relative genetic distance between two linked genes.



Haplotype

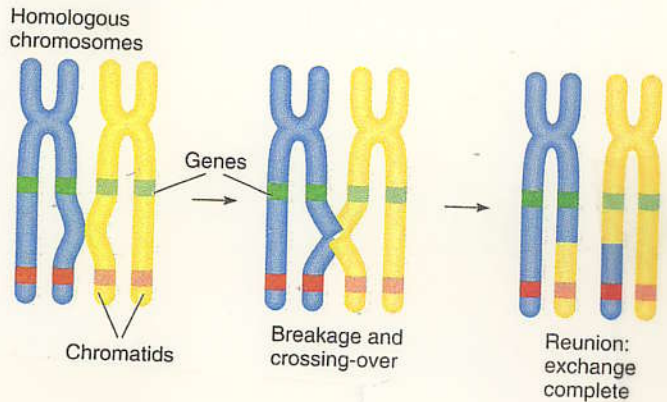
AB
vs.
ab

NO CROSSING OVER
< 5kb

∴
HAPLOTYPE
OR
COMPLEX
POLYMORPHIC
LOCUS

Figure 5.2

Mechanism of crossing-over. A highly simplified diagram of a crossover between two nonsister chromatids during meiotic prophase, giving rise to recombinant (nonparental) combinations of linked genes.



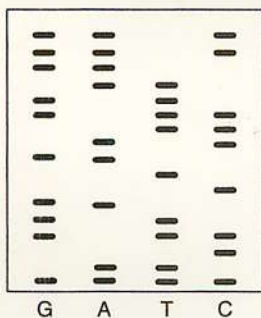
A Haplotype is a closely linked set of specific SNPs



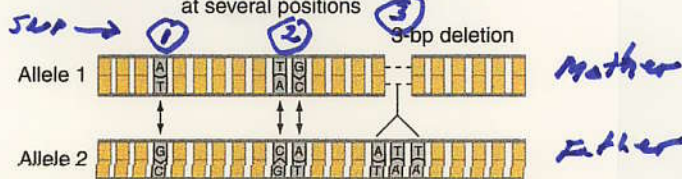
PCR amplification of HLA locus from one person who is heterozygous for two complex haplotypes

Clone from PCR products

Sequence several clones to obtain at least one sequence from each of the two alleles.



Production of two classes of clones that differ at several positions



3 SNP differences on each chromosome

They are Always inherited Together & Reflect Ancestry!

Haplotype

① ② -
①' ②' ③

Figure 9.15 The variations associated with a complex haplotype are best defined by sequencing. Using automated protocols to sequence an entire polymorphic region is often the most rapid and accurate way to detect changes associated with polymorphic alleles at a complex locus.

Mark Chromosome History - inherited as a unit

HAPLOTYPES

TRACE ORIGINS OF HUMANS

IMMORTAL Genetic VARIATION!

Complex Haplotypes

A contraction of the phrase "haploid genotype," the term **haplotype** refers to a specific combination of linked alleles in a cluster of related genes. Immunogeneticists often use it to describe the combination of alleles of the *major histocompatibility complex (MHC)*: a large cluster of genes on human chromosome 6 that play a role in the immune response. With the resolving power to look at DNA at the level of nucleotides, "haplotype" now refers to any set of linked DNA changes along a chromosome. These changes could be in one or several genes, or in noncoding stretches. The **complex** refers to the multiple types of variation that can exist at alternative alleles, including more than one nucleotide substitution, a substitution in combination with a small deletion, duplication, or other insertion. Thus, a **complex haplotype** is a set of linked DNA variations along a chromosome, with the possibility of many differences between alternative alleles.

H1
H2
H1
H2
H1

HAPLOTYPE PATTERNS	
Person A	ATTGAT CCGGAT...CCATCGGA...CTAA
Person B	ATTGAT AGGAT...CCAGCGGA...CTCA
Person C	ATTGAT CCGGAT...CCATCGGA...CTAA
Person D	ATTGAT AGGAT...CCAGCGGA...CTCA
Person E	ATTGAT CCGGAT...CCATCGGA...CTAA

Closely Linked SNPs

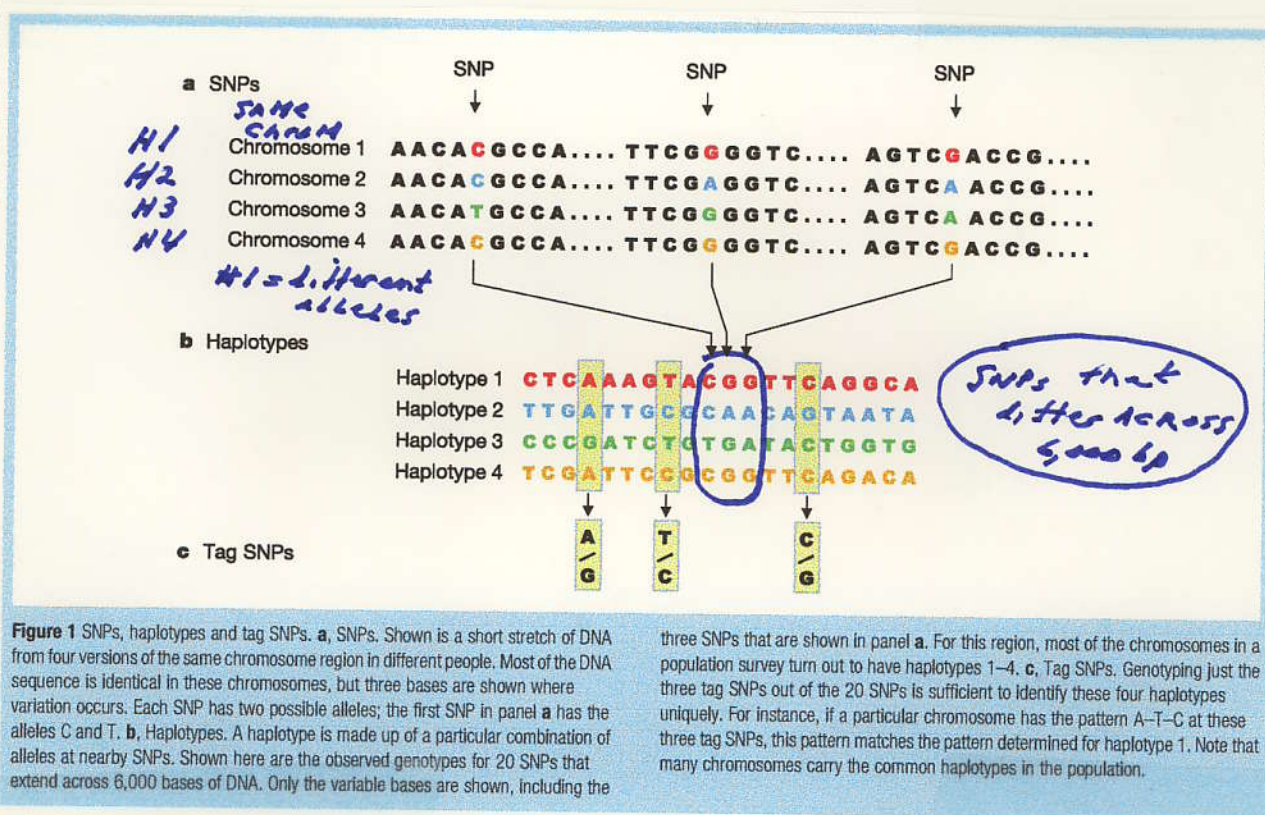
Building blocks. Persons B and D share a haplotype unlike the other three, characterized by three different SNPs.

The International HapMap Project

The International HapMap Consortium*

*Lists of participants and affiliations appear at the end of the paper

The goal of the International HapMap Project is to determine the common patterns of DNA sequence variation in the human genome and to make this information freely available in the public domain. An international consortium is developing a map of these patterns across the genome by determining the genotypes of one million or more sequence variants, their frequencies and the degree of association between them, in DNA samples from populations with ancestry from parts of Africa, Asia and Europe. The HapMap will allow the discovery of sequence variants that affect common disease, will facilitate development of diagnostic tools, and will enhance our ability to choose targets for therapeutic intervention.



THE 0.1% THAT'S DIFFERENT!

CORRELATE WITH SEQUENCE VARIANTS AFFECTING DISEASE

40

International HapMap Groups

news feature



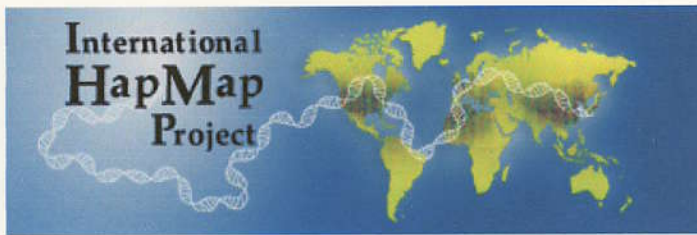
RIGHT: I. SLATER/CORBIS; FAR RIGHT: G. FRANKEN/CORBIS



World view: the HapMap initiative will gather genetic data from African, Asian and ancestrally European populations.

- ① Northern European
- ② Western European
- ③ Yoruba/Nigerian/African
- ④ Japanese & Han Chinese/Asian

"GROUP" Genetic Diversity to Disease & Other Aspects of Biology.



International HapMap Project

Home | About the Project | Data | Publications | Tutorial

中文 | [English](#) | Français | 日本語 | Yoruba

About the HapMap

[What is the HapMap?](#)

[Origins of Haplotypes](#)

[Health Benefits](#)

[Populations Sampled](#)

[Ethical Issues](#)

[Consent Forms](#)

[Community Advisory Groups\(CAG\)](#)

[Data Release Policy](#)

[Guidelines For Data Use](#)

[Guidelines For Referring to HapMap Populations](#)

Project Information

[About the Project](#)

[HapMap Publications](#)

[HapMap Tutorial](#)

[HapMap Mailing List](#)

[HapMap Project Participants](#)

[HapMap Mirror Site in Japan](#)

Useful Links

[HapMap Project Press Release](#)

[NHGRI HapMap Page](#)

[NCBI Variation Database \(dbSNP\)](#)

[Japanese SNP Database \(JSNP\)](#)

What Is the HapMap?

The HapMap is a catalog of common genetic variants that occur in human beings. It describes what these variants are, where they occur in our DNA, and how they are distributed among people within populations and among populations in different parts of the world. The International HapMap Project is not using the information in the HapMap to establish connections between particular genetic variants and diseases. Rather, the Project is designed to provide information that other researchers can use to link genetic variants to the risk for specific illnesses, which will lead to new methods of preventing, diagnosing, and treating disease.

The DNA in our cells

```

GAATAAATAAGTITTCCTGCTTCCTATTTGTCGTTACTTCGAATTTATTTATTTAATATTTATTTT
AGACGGAGTTTCACTCTGTGGCCAACTGGAGTCAAGTGGGTGATGTCAGCTCACTGCACACTCGCTTCTGG
TTCAAGCEATTTCTCTGCTCAGCCCTCTGAGTACGTTGGACTACAGTCACACACCACCGCCGCGCTTCTGG
TAITTTAGTAGAGTGGGGTTCCACATGTGGCCAGACTGGTTCGAACTCTGATCGATCAATTTGATCGCCAGCCT
GCCTCCCAADAGCTGGGATTACGGCCTGAGCCACCGCCGCTGGCCCTTTCGATCAATTTGATCGCTTCTTCT
TGCCCTGGACTTTACAAGCTTACCTGTCTGCTTCAGATTTTGTGGTCTCACTGCTGGCCAGTACGTAAGAAA
ATCCATGATTTGCTCATCCACTCCTGTTGTCATCTCCTTATCTGGGGTCCACTGCTCTCTCTGATGGATT
CTGATCCCCAGTACTAGCATGGCGTAACAACCTGGCTCTGCTTCCAGGGTCAAGCTGGGGTGGCTGTTCAAGC
TCAGAAAAATGCCATGTAAATTAATAAGATTTAAATATAGGAAAAAAGTAAGCAACAAGGAACAAAA
GGAAAGAACATGTATTCATCCATTATTTATATACAATTAAGAAATTTGGAACTTTAGATTACACTGCTTTAGAG
ATGGAGTGTAGTAACTTTTACTCTTACAAAATACATGTGTAGGAAATTTGGGAAGAAATAGTAACCTACCCAAA
CAGTGTAAATGTGAATATGTCCTACTAGAGGAAAGAGGCACTTGAAGAACATCTGTAACCCTATAAAGCAATTA
CATCATAATGTAAGAAACCAAGGAAATTTTATAGAAACATACAGGCTAAATAACAAGTAGAGCCACTGTCAAT
TTATCTCCCTTTGTGCTGTGTGAGAACTTAGAGTATATTTGACATAGCATGGAAAAATGGAGGCTAGTTATC
AACAGTTCATTTTAAAGTCAACACATCTAGTATAGGTGAAGTGTCTCTGCCAATGATTTGCCACTTTGTGC
CCAGATCCAGCATAGGGTATGTTGCCATTTACAACGTTTATGCTTAAAGAGGAAATTAAGAGCAAAACAGT
GCATGCTGGAGAGAAAGCTGATACAAATATAAATGAAACAATAATGGAAAAATTTGAGAAACTACTATTTCTAA
ATTAATCATGATTTTCTAGAAATTAAGCTTTTAAATTTTATGATAAATCCCAATGTGAGACAGATAAGTATTAGTAT
GGTATGAGTAAATAATCTGTATATAAATATTCATTTTCATAGTGGAAAGAAATAAATAAGGTTGTGATGTTGTG
ATATTTTCTAGAGGGTGTGAGGAAAGAAATGCTTTTTCATCTCTCTTCCACTAAGAAAGTTCGAACTATTT
AATTTAGGCACATCAATAATTAAGTCACTTAAATGCGAAAGGTAATTAAGAGACTTAAAGCTGAAAGTTA
AGTAGTGCACACTGAACTATATAAATAATGACAGGGTGGTGGAACTAGGCTTATATTAAGAGGGCTAAAAATG
CAATAAGACCACAGGCTTAAATGCTTTAAACTGTGAAAGGTGAAACTAGAAATGAATAAATCTATAAATTTG
ATCAAAAGAAAGAAACAAATTAAGTAAATTAATACAAGAAATGTTGGCTGGATCTGAGCAATGATTAAGT
AAGATAAACAAGAAATTTCTGAAATCTGGAAATCTTTGGGCTAACCTGAAACAGTATTTGAAACTATTT
TAAATGACAGTATAGTAAATTTTGAATCTATATGA
    
```



Figure 1: When DNA sequences on a part of chromosome 7 from two random individuals are compared, two single nucleotide polymorphisms (SNPs) occur in about 2,200 nucleotides.

contains long chains of four chemical building blocks -- adenine, thymine, cytosine, and guanine, abbreviated A, T, C, and G. More than 6 billion of these chemical bases, strung together in 23 pairs of chromosomes, exist in a human cell. (See <http://www.dnafb.org/dnafb/> for basic information about genetics.) These genetic sequences contain information that influences our physical traits, our likelihood of suffering from disease, and the responses of our bodies to substances that we encounter in the environment.

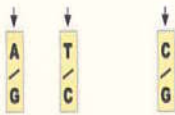


Figure 2: The construction of the HapMap occurs in three steps. (a) Single nucleotide polymorphisms (SNPs) are identified in DNA samples from multiple individuals. (b) Adjacent SNPs that are inherited together are compiled into "haplotypes." (c) "Tag" SNPs within haplotypes are identified that uniquely identify those haplotypes. By genotyping the three tag SNPs shown in this figure, researchers can identify which of the four haplotypes shown here are present in each individual.

The genetic sequences of different people are remarkably similar. When the chromosomes of two humans are compared, their DNA sequences can be identical for hundreds of bases. But at about one in every 1,200 bases, on average, the sequences will differ (Figure 1). One person might have an A at that location, while another person has a G, or a person might have extra bases at a given location or a missing segment of DNA. Each distinct "spelling" of a chromosomal region is called an allele, and a collection of alleles in a person's chromosomes is known as a genotype.

Differences in individual bases are by far the most common type of genetic variation. These genetic differences are known as single nucleotide polymorphisms, or SNPs (pronounced "snips"). By identifying most of the approximately 10 million SNPs estimated to occur commonly in the human genome, the International HapMap Project is identifying the basis for a large fraction of the genetic diversity in the human species.

For geneticists, SNPs act as markers to locate genes in DNA sequences. Say that a spelling change in a gene increases the risk of suffering from high blood pressure, but researchers do not know where in our chromosomes that gene is located. They could compare the SNPs in people who have high blood pressure with the SNPs of people who do not. If a particular SNP is more common among people with hypertension, that SNP could be used as a pointer to locate and identify the gene involved in the disease.

However, testing all of the 10 million common SNPs in a person's chromosomes would be extremely expensive. The development of the HapMap will enable geneticists to take advantage of how SNPs and other genetic variants are organized on chromosomes. Genetic variants that are near each other tend to be inherited together. For example, all of the people who have an A rather than a G at a particular location in a chromosome can have identical genetic variants at other SNPs in the chromosomal region surrounding the A. These regions of linked variants are known as haplotypes (Figure 2).

In many parts of our chromosomes, just a handful of haplotypes are found in humans. [See [The Origins of Haplotypes](#).] In a given population, 55 percent of people may have one version of a haplotype, 30 percent may have another, 8 percent may have a third, and the rest may have a variety of less common haplotypes. The International HapMap Project is identifying these common haplotypes in four populations from different parts of the world. It also is identifying "tag" SNPs that uniquely identify these haplotypes. By testing an individual's tag SNPs (a process known as genotyping), researchers will be able to identify the collection of haplotypes in a person's DNA. The number of tag SNPs that contain most of the information about the patterns of genetic variation is estimated to be about 300,000 to 600,000, which is far fewer than the 10 million common SNPs.

Once the information on tag SNPs from the HapMap is available, researchers will be able to use them to locate genes involved in medically important traits. Consider the researcher trying to find genetic variants associated with high blood

pressure. Instead of determining the identity of all SNPs in a person's DNA, the researcher would genotype a much smaller number of tag SNPs to determine the collection of haplotypes present in each subject. The researcher could focus on specific candidate genes that may be associated with a disease, or even look across the entire genome to find chromosomal regions that may be associated with a disease. If people with high blood pressure tend to share a particular haplotype, variants contributing to the disease might be somewhere within or near that haplotype.

[Home](#) | [About the Project](#) | [Data](#) | [Publications](#) | [Tutorial](#)

Please send questions and comments on website to
help@hapmap.org



International HapMap Project

[Home](#) | [About the Project](#) | [Data](#) | [Publications](#) | [Tutorial](#)

[中文](#) | [English](#) | [Français](#) | [日本語](#) | [Yoruba](#)

About the HapMap

[What is the HapMap?](#)

[Origins of Haplotypes](#)

[Health Benefits](#)

[Populations Sampled](#)

[Ethical Issues](#)

[Consent Forms](#)

[Community Advisory Groups\(CAG\)](#)

[Data Release Policy](#)

[Guidelines For Data Use](#)

[Guidelines For Referring to HapMap Populations](#)

Project Information

[About the Project](#)

[HapMap Publications](#)

[HapMap Tutorial](#)

[HapMap Mailing List](#)

[HapMap Project Participants](#)

[HapMap Mirror Site in Japan](#)

Useful Links

[HapMap Project Press Release](#)

[NHGRI HapMap Page](#)

[NCBI Variation Database \(dbSNP\)](#)

[Japanese SNP Database \(JSNP\)](#)

Which Populations Are Being Sampled

The International HapMap Project is analyzing DNA from populations with African, Asian, and European ancestry. Together, these DNA samples should enable HapMap researchers to identify most of the common haplotypes that exist in populations worldwide. [[See What Is the HapMap?](#)]

Because of the history of the human species, most of the common haplotypes in human chromosomes occur in all human populations. [[See The Origin of Haplotypes.](#)] However, any given haplotype may be more common in one population and less common in another, and newer haplotypes may be found in just a single population. Efficiently choosing the tag SNPs needed to identify haplotypes therefore requires looking at haplotype frequencies in multiple populations. Also, genetic data from more than one population will enhance the ability of researchers to study the genetic contributions to diseases that are more or less prevalent in different groups.

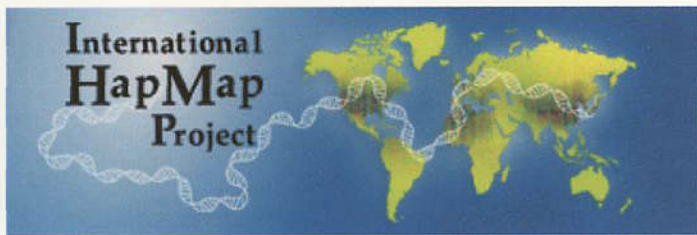
The DNA samples for the HapMap have come from a total of 270 people. The Yoruba people of Ibadan, Nigeria, provided 30 sets of samples from two parents and an adult child (each such set is called a trio). In Japan, 45 unrelated individuals from the Tokyo area provided samples. In China, 45 unrelated individuals from Beijing provided samples. Thirty U.S. trios provided samples, which were collected in 1980 from U.S. residents with northern and western European ancestry by the Centre d'Etude du Polymorphisme Humain (CEPH).

The blood samples are being converted into cell lines, which are used to make DNA, by the non-profit Coriell Institute for Medical Research. [[See <http://locus.umdj.edu/nigms> for more information.](#)] Coriell provides DNA and cell lines from the samples for research projects that have been approved by the appropriate ethics committees. The samples and cell lines are not linked to any individual in the populations studied. However, the samples and cell lines are identified as coming from one of the four populations participating in the study, which raises ethical issues associated with conducting genetic research in named populations. [[See \[How Are Ethical Issues Being Addressed?\]\(#\) and \[Guidelines for Referring to the HapMap Populations in Publications and Presentations.\]\(#\)](#)]

To assess how much additional information would be gained by genotyping other populations, haplotypes in a set of chromosomal regions are being analyzed in samples from several additional populations.

[Home](#) | [About the Project](#) | [Data](#) | [Publications](#) | [Tutorial](#)

Please send questions and comments on website to help@hapmap.org



International HapMap Project

[Home](#) | [About the Project](#) | [Data](#) | [Publications](#) | [Tutorial](#)

[中文](#) | [English](#) | [Français](#) | [日本語](#) | [Yoruba](#)

About the HapMap

[What is the HapMap?](#)

[Origins of Haplotypes](#)

[Health Benefits](#)

[Populations Sampled](#)

[Ethical Issues](#)

[Consent Forms](#)

[Community Advisory Groups\(CAG\)](#)

[Data Release Policy](#)

[Guidelines For Data Use](#)

[Guidelines For Referring to HapMap Populations](#)

Project Information

[About the Project](#)

[HapMap Publications](#)

[HapMap Tutorial](#)

[HapMap Mailing List](#)

[HapMap Project Participants](#)

[HapMap Mirror Site in Japan](#)

Useful Links

[HapMap Project Press Release](#)

[NHGRI HapMap Page](#)

[NCBI Variation Database \(dbSNP\)](#)

[Japanese SNP Database \(JSNP\)](#)

How Will the HapMap Benefit Human Health?

The International HapMap Project will benefit human health by providing an extensive resource that researchers can use to discover the genetic variants involved in disease and individual responses to therapeutic agents. Once such variants have been discovered, researchers can learn much more about the origins of illnesses and about ways to prevent, diagnose, and treat those illnesses.

The goal of the Project is not to identify these disease-related genes directly. Rather, by identifying haplotypes, the HapMap provides a tool that can be used in what are called association studies. For these studies, researchers will compare the haplotypes in individuals with a disease to the haplotypes of a comparable group of individuals without a disease (the controls). If a particular haplotype occurs more frequently in affected individuals compared with controls, a gene influencing the disease may be located within or near that haplotype.

Common diseases such as cancer, stroke, heart disease, diabetes, depression, and asthma usually result from the combined effects of a number of genetic variants and environmental factors. According to an idea known as the common disease-common variant hypothesis, the risk of contracting common diseases is influenced by genetic variants that are relatively common in populations. Not enough data are yet available to evaluate the generality of this hypothesis, but more and more widely distributed genetic variants associated with common diseases are being discovered, including variants that contribute to autoimmune diseases, schizophrenia, diabetes, asthma, stroke, and heart attacks. One of the many benefits of the International HapMap Project will be the use of the HapMap to learn more about the links between these common disorders and our genes.

Knowledge derived from use of the HapMap also will result in advances that are difficult to predict today. Medical treatments could be customized, based on a patient's genetic make-up, to maximize effectiveness and minimize side effects. Genetic variants contributing to longevity or resistance to disease could be identified, leading to new therapies with widespread benefits. As with any new body of knowledge, the HapMap is likely to lead to both new challenges and to unexpected and unprecedented opportunities.

[Home](#) | [About the Project](#) | [Data](#) | [Publications](#) | [Tutorial](#)

Please send questions and comments on website to help@hapmap.org

Human pigmentation genetics: the difference is only skin deep

Richard A. Sturm,^{1*} Neil F. Box,¹ and Michele Ramsay²

Summary

There is no doubt that visual impressions of body form and color are important in the interactions within and between human communities. Remarkably, it is the levels of just one chemically inert and stable visual pigment known as melanin that is responsible for producing all shades of humankind. Major human genes involved in its formation have been identified largely using a comparative genomics approach and through the molecular analysis of the pigmentary process that occurs within the melanocyte. Three classes of genes have been examined for their contribution to normal human color variation through the production of hypopigmented phenotypes or by genetic association with skin type and hair color. The MSH cell surface receptor and the melanosomal P-protein are the two most obvious candidate genes influencing variation in pigmentation phenotype, and may do so by regulating the levels and activities of the melanogenic enzymes tyrosinase, TRP-1 and TRP-2. *BioEssays* **20**:712-721, 1998. © 1998 John Wiley & Sons, Inc.

TABLE 1. Human Pigmentation Genes

Gene symbol	Mouse homologue	Chromosome	Phenotype	Protein	Function/activity
TYR	Albino (c)	11q14-21	OCA1	Tyrosinase	Tyrosine hydroxylation; DOPA oxidase
TYRP1	Brown (b)	9p23	OCA3/ROCA	TRP-1	DHICA oxidase
TYRP2	Slaty (slt)	13q31-32	Unknown	TRP-2	Dopachrome tautomerase
P	Pink-eyed dilute (p)	15q11.2-12	OCA2, BOCA	P-protein	Melanosomal transmembrane protein
MC1R	Extension (e)	16q24.3	Red hair	MSHR	G-protein-coupled receptor

(+)
JLQ24A5

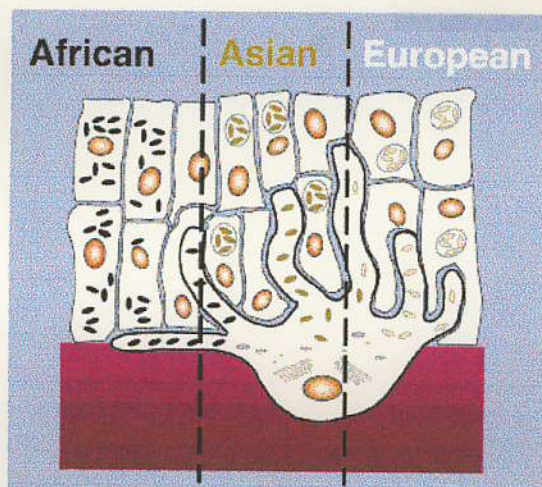


Figure 1. Variation in melanosome structure and distribution in different groups. A single skin melanocyte cell interdigitating with keratinocyte cells is partitioned into three sections. Shown within the melanocyte are the four stages of melanosome formation from budding from the Golgi apparatus, to the fully pigmented stage IV melanosomes migrating up the dendritic processes of the cell and secreted into the keratinocytes. In African populations, the melanosomes remain as singular heavily pigmented particles while in Asians and Europeans the melanosomes cluster in membrane bound organelles giving different skin complexions.

SLC24A5, a Putative Cation Exchanger, Affects Pigmentation in Zebrafish and Humans

Rebecca L. Lamason,^{1*} Manzoor-Ali P.K. Mohideen,^{1†} Jason R. Mest,¹ Andrew C. Wong,^{1‡} Heather L. Norton,⁶ Michele C. Aros,¹ Michael J. Jurynecek,⁸ Xianyun Mao,⁶ Vanessa R. Humphreville,^{1§} Jasper E. Humbert,^{2,9} Soniya Sinha,² Jessica L. Moore,^{1||} Pudur Jagadeeswaran,¹⁰ Wei Zhao,³ Gang Ning,⁷ Izabela Makalowska,⁷ Paul M. McKeigue,¹¹ David O'Donnell,¹¹ Rick Kittles,¹² Esteban J. Parra,¹³ Nancy J. Mangini,¹⁴ David J. Grunwald,⁸ Mark D. Shriver,⁶ Victor A. Canfield,⁴ Keith C. Cheng^{1,4,5¶}

Lighter variations of pigmentation in humans are associated with diminished number, size, and density of melanosomes, the pigmented organelles of melanocytes. Here we show that zebrafish *golden* mutants share these melanosomal changes and that *golden* encodes a putative cation exchanger *slc24a5* (*nckx5*) that localizes to an intracellular membrane, likely the melanosome or its precursor. The human ortholog is highly similar in sequence and functional in zebrafish. The evolutionarily conserved ancestral allele of a human coding polymorphism predominates in African and East Asian populations. In contrast, the variant allele is nearly fixed in European populations, is associated with a substantial reduction in regional heterozygosity, and correlates with lighter skin pigmentation in admixed populations, suggesting a key role for the *SLC24A5* gene in human pigmentation.

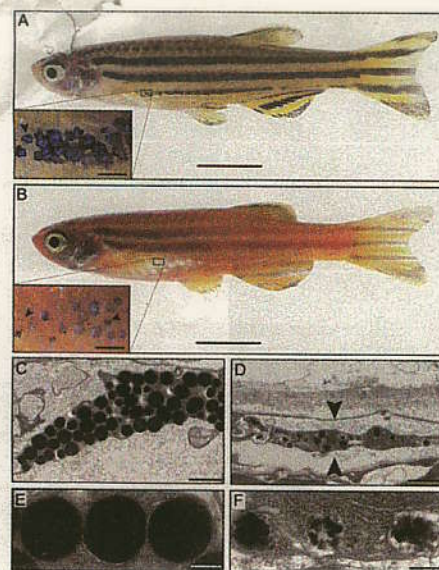


Fig. 1. Phenotype of *golden* zebrafish. Lateral views of adult wild-type (A) and *golden* (B) zebrafish. Insets show melanophores (arrowheads). Scale bars, 5 mm (inset, 0.5 mm). *gol^{b1}* mutants have melanophores that are, on average, smaller, more pale, and transparent. Transmission electron micrographs of skin melanophore from 55-hpf wild-type (C and E) and *gol^{b1}* (D and F) larvae. *gol^{b1}* skin melanophores (arrowheads show edges) are thinner and contain fewer melanosomes than do those of wild type. Melanosomes of *gol^{b1}* larvae are fewer in number, smaller, less-pigmented, and irregular compared with wild type. Scale bars in (C) and (D), 1000 nm; in (E) and (F), 200 nm.

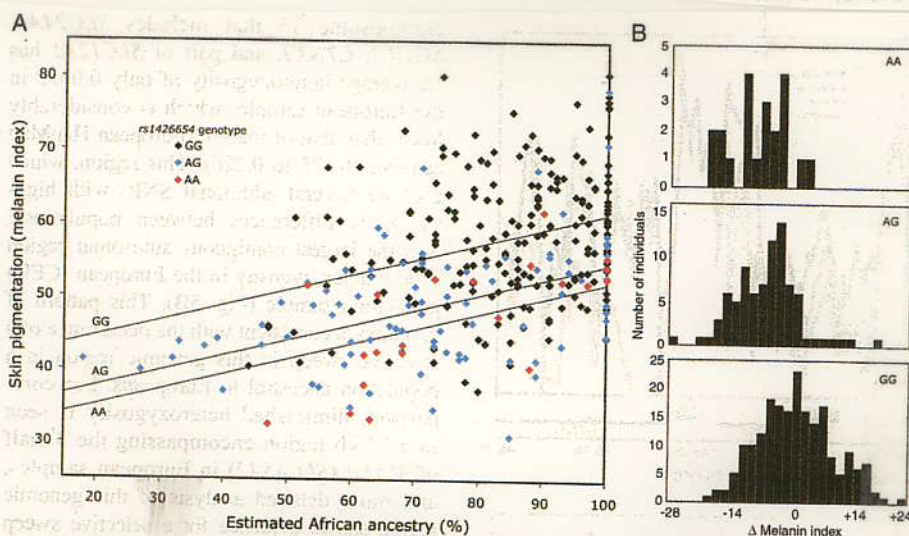


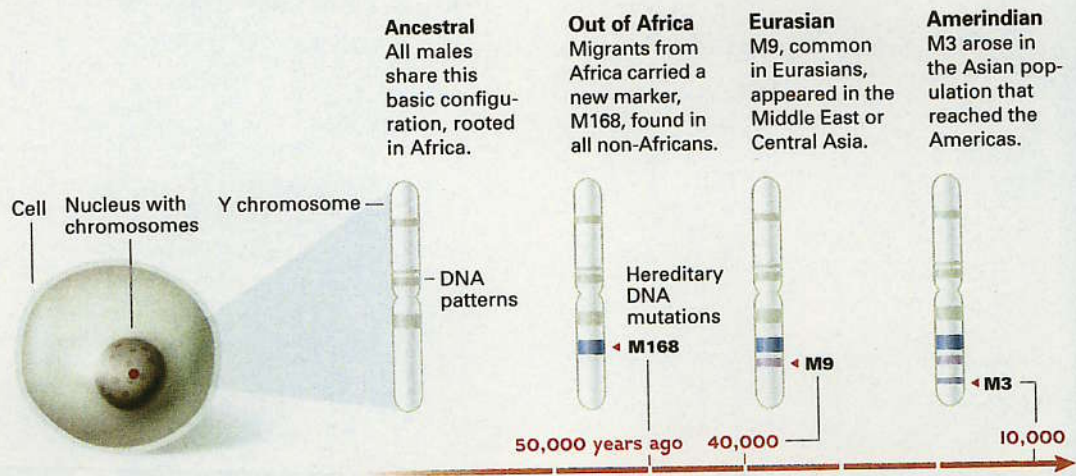
Fig. 6. Effect of *SLC24A5* genotype on pigmentation in admixed populations. (A) Variation of measured pigmentation with estimated ancestry and *SLC24A5* genotype. Each point represents a single individual; *SLC24A5* genotypes are indicated by color. Lines show regressions, constrained to have equal slopes, for each of the three genotypes. (B) Histograms showing the distribution of pigmentation after adjustment for ancestry for each genotype. Values shown are the difference between the measured melanin index and the calculated GG regression line ($y = 0.2113x + 30.91$). The corresponding uncorrected histograms are shown in fig. S7. Mean and SD (in parentheses) are given as follows: for GG, 0 (8.5), $n = 202$ individuals; for AG, -7.0 (7.4), $n = 85$; for AA, -9.6 (6.4), $n = 21$.



HAPLOTYPES CAN BE USED
TO TRACE ORIGINS
OF HUMAN POPULATIONS

History on a Chromosome

Genetic mutations act as markers, tracing a journey through time. The earliest known mutation to spread outside Africa is M168, which arose some 50,000 years ago. This graphic shows the Y chromosome of a Native American man with various mutations including M168, proving his African ancestry.



68 NATIONAL GEOGRAPHIC • MARCH 2006

NGM ART

Note - looking at novel alleles in populations! And their frequencies!

e.g. A vs. a

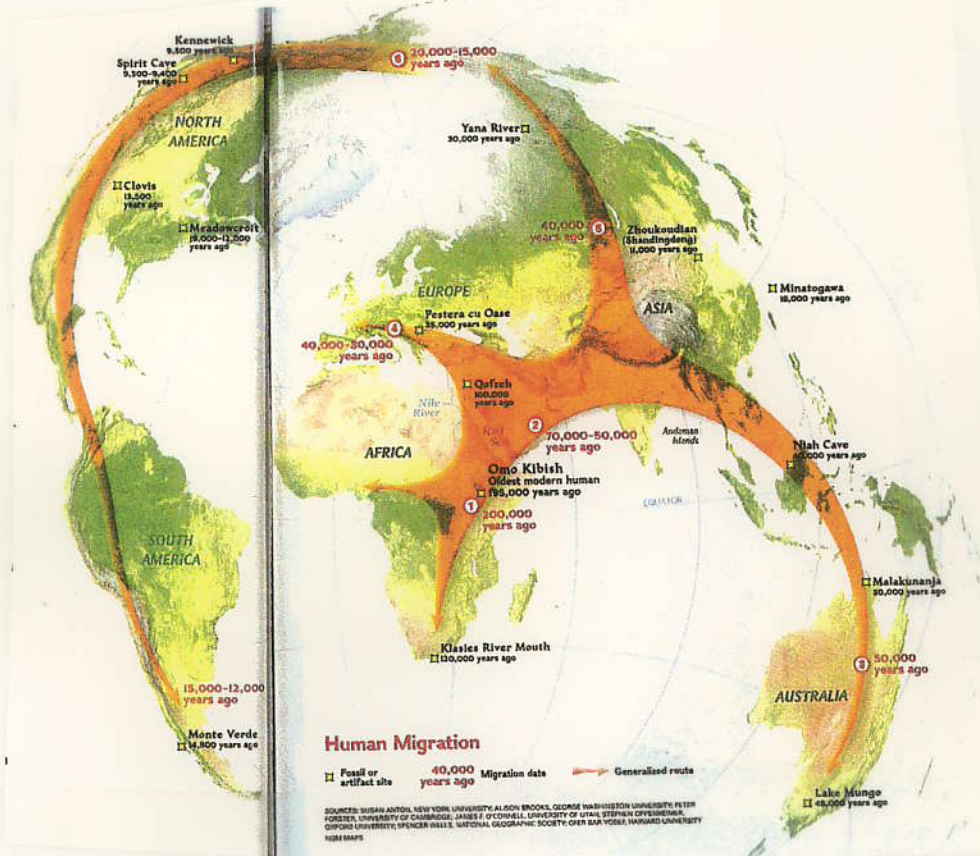
B vs b

C vs c

Markers MUST have different "forms"!

(52)

What Genes and Fossils Tell Us



1 African Cradle

Most paleoanthropologists and geneticists agree that modern humans arose some 200,000 years ago in Africa. The earliest modern human fossils were found in Omo Kibish, Ethiopia. Sites in Israel hold the earliest evidence of modern humans outside Africa, but that group went no farther, dying out about 90,000 years ago.

2 Out of Africa

Genetic data show that a small group of modern humans left Africa for good 70,000 to 50,000 years ago and eventually replaced all earlier types of humans, such as Neandertals. All non-Africans are the descendants of these travelers, who may have migrated around the top of the Red Sea or across its narrow southern opening.

3 The First Australians

Discoveries at two ancient sites—artifacts from Malakunanja and fossils from Lake Mungo—indicate that modern humans followed a coastal route along southern Asia and reached Australia nearly 50,000 years ago. Their descendants, Australian Aborigines, remained genetically isolated on that island continent until recently.

4 Early Europeans

Paleoanthropologists long thought that the peopling of Europe followed a route from North Africa through the Levant. But genetic data show that the DNA of today's western Eurasians resembles that of people in India. It's possible that an inland migration from Asia seeded Europe between 40,000 and 30,000 years ago.

5 Populating Asia

Around 40,000 years ago, humans pushed into Central Asia and arrived on the grassy steppes north of the Himalaya. At the same time, they traveled through Southeast Asia and China, eventually reaching Japan and Siberia. Genetic clues indicate that humans in northern Asia eventually migrated to the Americas.

6 Into the New World

Exactly when the first people arrived in the Americas is still hotly debated. Genetic evidence suggests it was between 20,000 and 15,000 years ago, when sea levels were low and land connected Siberia to Alaska. Ice sheets would have covered the interior of North America, forcing the new arrivals to travel down the west coast.

WATCH "BLACKBEARD" ON NG CHANNEL, SUNDAY, MARCH 12, 8 P.M. ET/9 P.M. PT

NATIONALGEOGRAPHIC.COM/MAGAZINE

MARCH 2006

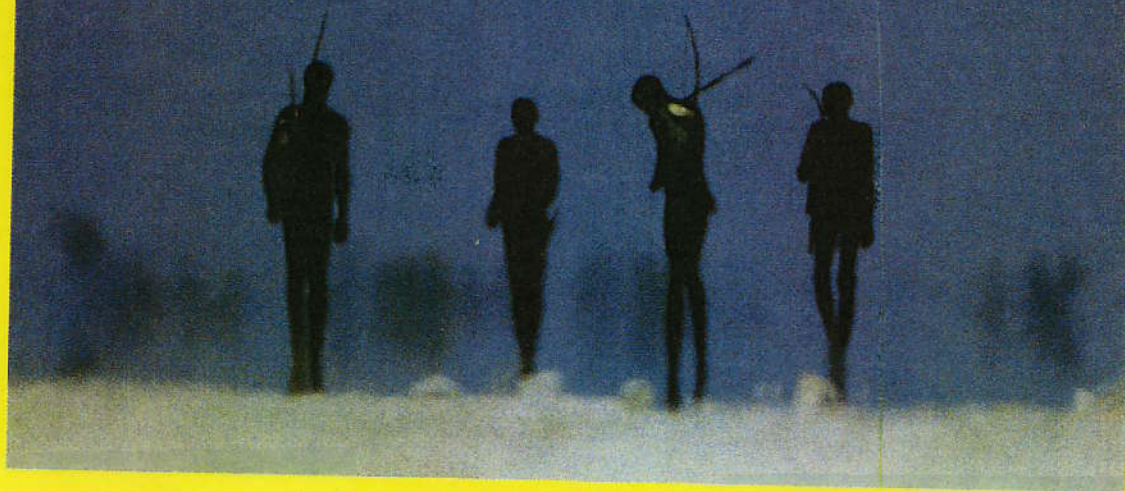
NATIONAL GEOGRAPHIC

The Greatest Journey Ever Told THE TRAIL OF OUR DNA

Ukraine's Revolution 32 Celtic Realm 74 The High Cost of Cheap Coal 96

Africa's Last Wolves 124 Battle of Hampton Roads 136

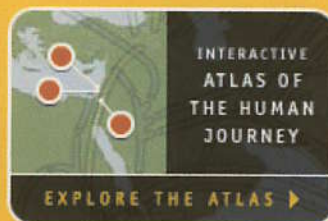
ZIP USA: Survival of the Richest 148



54

A LANDMARK STUDY OF THE HUMAN JOURNEY

Who was **your** first ancestor? New DNA studies say that all humans descended from an African ancestor who lived only 60,000 years ago. Uncover the specific paths that led from him to you—the ultimate human history, as written in our genes.



NEWS AND RESOURCES

- **Millions of Men May Be Descended From Irish King, Study Says**
- **More Related News**
- **Related Web Resources**

GENETICS OVERVIEW



Delve inside and explore the basics of genetics, from chromosomes to natural selection.

YOUR GENETIC JOURNEY



Explore your own genetic journey with Dr. Spencer Wells. DNA analysis includes a depiction of your ancient

ancestors and an interactive map tracing your genetic lineage around the world and through the ages.

Question of the Week

- **What happens once I order a Genographic Project Participation Kit?**

Enter Your Genographic Project Kit ID Here

LOG IN

- Remember my log in (optional)

ALSO SEE

- **Buy the Participation Kit**
- **Help Support the Genographic Project Field Research**



Global field offices supported by the Waitt Family Foundation



A research partnership of National Geographic and IBM

The Genographic Project

Published 2 weeks, 1 day ago

Two months ago my father told me that he had just sent a cheek swab sample to the [National Geographic Genographic Project](#), which, at the time, I had never heard of. To have the website tell it, the goal of the project is to “understand the human journey — where we came from and how we got to where we live today.” The data collected “will map world migratory patterns dating back some 150,000 years and will fill in the huge gaps in our knowledge of humankind’s migratory history.”

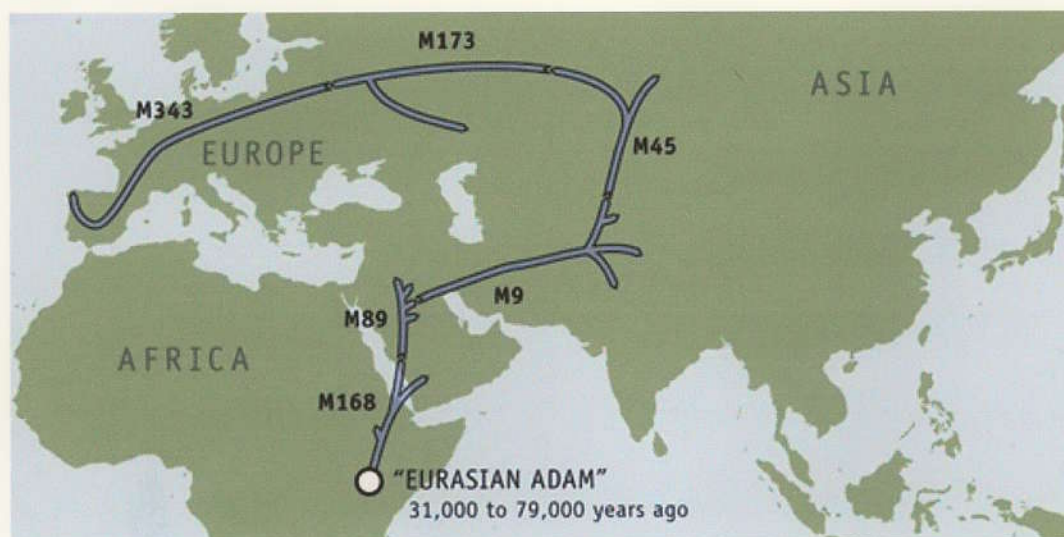
The [participation kit](#) is only \$100 and it’s money well spent if you ask me.

After having read both Jared Diamond’s [Guns, Germs, and Steel: The Fates of Human Societies](#) and Steve Olson’s [Mapping Human History: Genes, Race, and Our Common Origins](#), I couldn’t wait for my father to receive the results — I’m completely fascinated by this stuff.

The results

To be clear—these tests are not conventional genealogy. Your results will not provide names for your personal family tree or tell you where your great grandparents lived. Rather, they will indicate the maternal or paternal genetic markers your deep ancestors passed on to you and the story that goes with those markers.

A few days ago he got the results and shared the online account with me. They’re utterly fascinating. The picture below shows my “ancestral journey.” At the website, this picture is interactive — clicking on a marker will tell you the story behind it.



In addition to this migratory map, you’re also given your genetic sequence (not reproduced here) and told how to interpret it in light of these mass migratory patterns. Finally, you’re told your “genetic history” (reproduced below), which basically summarizes where and when your haplogroup originated, how they got

INTRODUCTION

ABOUT THE PROJECT

HOW TO PARTICIPATE

FREQUENTLY ASKED QUESTIONS

HOW TO PARTICIPATE

Public participation, including yours, is critical to the Genographic Project's success. Here's how you can get involved:

Purchasing a Public Participation Kit will fund important research around the world—and open the door to the ancient past of your own genetic background.

With a simple and painless cheek swab you can sample your own DNA. You'll submit the sample through our secure, private, and completely anonymous system, then log on to the project Web site to track your personal results online.

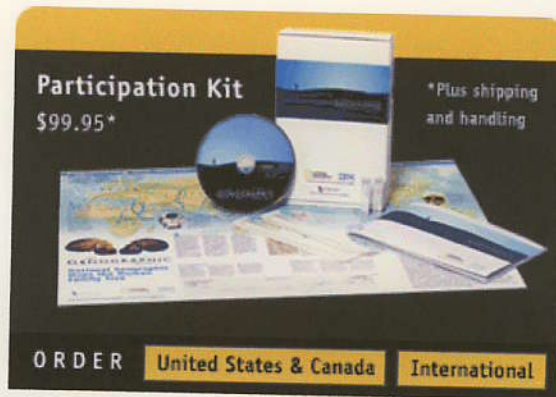
This is not a genealogy test and you won't learn about your great grandparents. You will learn, however, of your deep ancestry, the ancient genetic journeys and physical travels of your distant relatives.

To insure total anonymity you will be identified at all times only by your kit number, not by your name. There is no record, no database that links test results with the names of their contributors. If you lose the kit number there will be no way to access your genetic results.

As your own genetic ancestry is revealed you'll also see worldwide samples map humankind's shared genetic background around the world and through the ages.

If you'd like to contribute your own results to the project's global database you'll be asked to answer a dozen "phenotyping" questions that will help place your DNA in cultural context.

This process is optional and completely anonymous, but it's also important. Each of us has a part



The kit includes the following elements:

- Buccal swab kit
- Multimedia DVD
- Exclusive National Geographic Genographic Map
- "Quick Start" card
- Genographic Project Brochure
- Self-addressed envelope
- Confidential Genographic Project ID Number (GPID)

The purchase price also includes the cost of the testing and analysis. International participants please [see note below](#).

Jewish and Middle Eastern non-Jewish populations share a common pool of Y-chromosome biallelic haplotypes

M. F. Hammer^{*†‡}, A. J. Redd^{**}, E. T. Wood^{**}, M. R. Bonner^{*}, H. Jarjanazi^{*}, T. Karafet^{*}, S. Santachiara-Benerecetti[¶], A. Oppenheim^{||}, M. A. Jobling^{**}, T. Jenkins^{††}, H. Ostrer^{‡‡}, and B. Bonn -Tamir[§]

^{*}Laboratory of Molecular Systematics and Evolution, University of Arizona, Tucson, AZ 85721; [¶]Department of Genetics, Universit  degli Studi di Pavia, Pavia 27100, Italy; ^{||}Hadassah Medical School, Hebrew University of Jerusalem, Jerusalem 91120, Israel; ^{**}Department of Genetics, University of Leicester, Leicester LE1 7RH, England; ^{††}SAMIR, University of Witwatersrand, Johannesburg 2000, South Africa; ^{‡‡}Department of Pediatrics, New York University Medical Center, New York, NY 10016; and [§]Department of Human Genetics, Sackler School of Medicine, Ramat Aviv 69978, Israel

Communicated by Arno G. Motulsky, University of Washington, Seattle, WA, March 15, 2000 (received for review November 17, 1999)

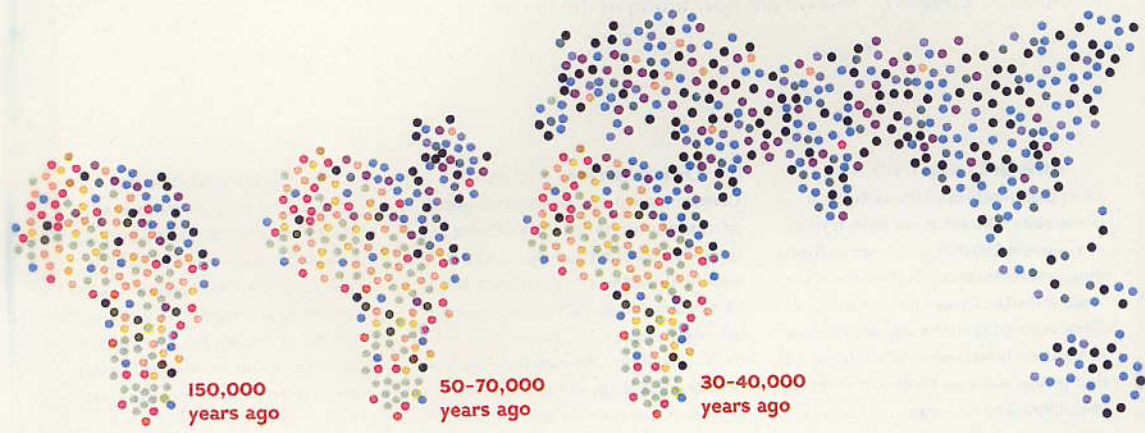
Haplotypes constructed from Y-chromosome markers were used to trace the paternal origins of the Jewish Diaspora. A set of 18 biallelic polymorphisms was genotyped in 1,371 males from 29 populations, including 7 Jewish (Ashkenazi, Roman, North African, Kurdish, Near Eastern, Yemenite, and Ethiopian) and 16 non-Jewish groups from similar geographic locations. The Jewish populations were characterized by a diverse set of 13 haplotypes that were also present in non-Jewish populations from Africa, Asia, and Europe. A series of analyses was performed to address whether modern Jewish Y-chromosome diversity derives mainly from a common Middle Eastern source population or from admixture with neighboring non-Jewish populations during and after the Diaspora. Despite their long-term residence in different countries and isolation from one another, most Jewish populations were not significantly different from one another at the genetic level. Admixture estimates suggested low levels of European Y-chromosome gene flow into Ashkenazi and Roman Jewish communities. A multidimensional scaling plot placed six of the seven Jewish populations in a relatively tight cluster that was interspersed with Middle Eastern non-Jewish populations, including Palestinians and Syrians. Pairwise differentiation tests further indicated that these Jewish and Middle Eastern non-Jewish populations were not statistically different. The results support the hypothesis that the paternal gene pools of Jewish communities from Europe, North Africa, and the Middle East descended from a common Middle Eastern ancestral population, and suggest that most Jewish communities have remained relatively isolated from neighboring non-Jewish communities during and after the Diaspora.



MOST GENETIC DIVERSITY
ORIGINATED IN THE
FOUNDER POPULATIONS
TO MODERN HUMANS!

Diverse From the Start

The diversity of genetic markers is greatest in Africa (multicolored dots in map), indicating it was the earliest home of modern humans. Only a handful of people, carrying a few of the markers, walked out of Africa (center) and, over tens of thousands of years, seeded other lands (right). "The genetic makeup of the rest of the world is a subset of what's in Africa," says Yale geneticist Kenneth Kidd.



© KENNETH K. KIDD

ARE There HUMAN
"Races?"

What is the History
& the Biology?



Race is Primarily a Sociological Concept that has caused much Human Suffering



Race is a largely non-biological concept con-founded by misunderstanding and a long history of prejudice. The relationship of genomics to the concepts of race and ethnicity has to be considered within complex historical and social contexts.

Most variation in the genome is shared between all populations, but certain alleles are more frequent in some populations than in others, largely as a result of history and geography. Use of genetic data to define racial groups, or of racial categories to classify biological traits, is prone to misinterpretation. To minimize such misinterpretation, the biological and sociocultural factors that inter-relate genetics with constructs of race and ethnicity need to be better understood and communicated within the next few years.

This will require research on how different individuals and cultures conceive of race, ethnicity, group identity and self-identity, and what role they believe genes or other biological factors have. It will also require a critical examination of how the scientific community understands and uses these concepts in designing research and presenting findings, and of how the media report these. Also necessary is widespread education about the biological meaning and limitations of research findings in this area (Box 6) and the formulation and adoption of public-policy options that protect against genomics-based discrimination or maltreatment (see Grand Challenge III-1).

Genome Project Goal

- ① MOST VARIATION SHARED
- ② Differences due to Migrations + Geographic Isolation
- ③ Most VARIATION SHARED by groups!!

Based on a very few genes that vary between groups FAR MORE than Majority - VAST Majority of other Genes

A genetic melting-pot

Marcus W. Feldman, Richard C. Lewontin and Mary-Claire King

Race as a biological concept has had a variety of meanings. In the taxonomic literature, a race is any distinguishable type within a species, such as dark-bellied and light-bellied variants of small mammals. In 1937, Theodosius Dobzhansky introduced the idea of geographical races — populations of species that differ in the frequencies of one or more genetic variants. But as no two populations have identical gene frequencies at variable (polymorphic) loci, Dobzhansky's definition of race becomes synonymous with that of population.

The classical definition of race, as applied to our species, is based on phenotypes such as skin colour, facial features and hair form that clearly differ between native inhabitants of different regions of the world. An underlying assumption is that all of these defining features (all largely genetic traits, although few of their genes have been identified) are characteristic of the genome in general. In other words, just as there are large differences between races in genes for skin colour, so there should be large genetic differences between races in general. In the previous absence of data to confirm or deny this assumption, it was not an unreasonable one to make.

But recent studies of genetic diversity indicate that the genes underlying the phenotypic differences used to assign race categories are atypical, in that they vary between races much more than genes in general. Together, the iconic features of race correlate well with continent of origin but do not reflect genome-wide differences between groups.

Discussion has arisen over the implications of these findings for the utility of racial classification in medical practice. The issue of whether race is a biologically useful or even meaningful concept when applied to humans in a medical context is controversial — holders of opposing views each claim to have evidence to support them. But there is no contradiction between these two well-substantiated bodies of data, as they actually



Human migration, as depicted in Charles Hunt's *Landing at Madras*, blurs the boundaries of race.

deal with two different questions that have become confused with one another.

The first question is: "Is it possible to find DNA sequences that differ sufficiently between populations to allow correct assignment of major geographical origin with high probability?" The answer to this question is yes, as shown by studies of genetic polymorphisms and by universal personal experience.

The second question is: "What fraction of human genetic variation, whether based on protein-coding genes or other sequences, falls within geographically separated populations, and what fraction occurs between these populations?" The answer to this question is that most genetic diversity occurs within groups, and that very little is found between them.

Why this apparent paradox? The answer is that genes that are geographically distinctive in their frequencies are not typical of the human genome in general.

It has been suggested that racial categorization has a valid role in good medical practice because many medically important genes vary between populations from different regions. But although knowing a patient's ancestry is often extremely useful in diagnosis and treatment, race is both too broad and too narrow a definition of ancestry to be biologically useful.

For any species, definitions of race can lose their discriminating power when individuals migrate to different regions and mate with their counterparts there. Among humans, large-scale migrations between continents — particularly through European colonial expansion and the commercial slave trade — has resulted in matings of individuals from different continents and the creation of new populations, especially in the Western Hemisphere and Oceania. Many people thus have ancestry from more than one major geographical region, meaning that the association of phenotype and geography breaks down.

For example, sickle-cell disease, which is often thought to be an African trait, is instead characteristic of ancient ancestry in a geographic region where malaria was endemic. Africa is one such region, but so are the Mediterranean and southern India. If sickle-cell disease is suspected, then the correct diagnostic approach is not simply to determine the patient's race, but to ask whether they have African, Mediterranean or South Indian ancestry. To use genotype effectively in making diagnostic and therapeutic decisions, it is not race that is relevant, but both intra- and trans-continental contributions to a person's ancestry.

Race and ancestry are confounded both by genetic heterogeneity within groups and by the widespread mixing of previously iso-

Race

Ancestral genetic data are far more useful for medical purposes than are racial categories, which may be correlated with disease for social or economic rather than biological reasons.

lated populations. The assignment of a racial classification to an individual hides the biological information that is needed for intelligent therapeutic and diagnostic decisions. A person classified as 'black' or 'Hispanic' by social convention could have any mixture of ancestries, as defined by continent of origin. Confusing race and ancestry could be potentially devastating for medical practice.

Other attempts to classify people into broad genetic groups based on the frequency of specific genes for, say, drug-metabolizing enzymes, are also likely to be poor predictors of medical outcome. As with racial groupings, the overall variation in the frequencies of such genes between groups is likely to be less than that within each group.

The conventional, social definition of race is useful in a medical context as it provides information about the social circumstances and lifestyle of patients. But this is a consequence of social history, so any variation is (at least in principle) transitory. By contrast, information on the likelihood that a person carries specific disease-related or treatment-response genes is grounded in their ancestry in far more complex ways. We suggest that identifying all contributions to a patient's ancestry can be useful in diagnosing and treating diseases with genetic influences. Eventually, for both diagnosis and treatment, specific genetic variants will provide concrete, useful information.

Marcus W. Feldman is in the Department of Biological Sciences, Stanford University, California 94305, USA.

Richard C. Lewontin is at the Museum of Comparative Zoology, Harvard University, Cambridge, Massachusetts 02138, USA.

Mary-Claire King is in the Departments of Genome Science and Medicine, University of Washington, Seattle, Washington 98195, USA.

FURTHER READING

- Rosenberg, N. A. *et al.* *Science* **298**, 2381–2385 (2002).
- González Burchard, E. *et al.* *N. Engl. J. Med.* **348**, 1170–1175 (2003).
- Cooper, R. S., Kaufman, J. S. & Ward, R. *N. Engl. J. Med.* **348**, 1166–1170 (2003).
- Lewontin, R. C. *Evol. Biol.* **6**, 381–398 (1972).
- Barbujani, G., Magagni, A., Minch, E. & Cavalli-Sforza, L. L. *Proc. Natl Acad. Sci. USA* **94**, 4516–4519 (1997).

HUMAN DIVERSITY

RICHARD LEWONTIN

*Scientific American Library
1992 ISBN 07167-1469-8*



ARE There
HUMAN RACES
or
Human BEINGS?

DNA testing CAN provide
information on Human
origins - geographical
Groupings + Ancestry!

WE WILL Learn all about
HUMAN VARIATION FROM GENOME SEQUENCING!

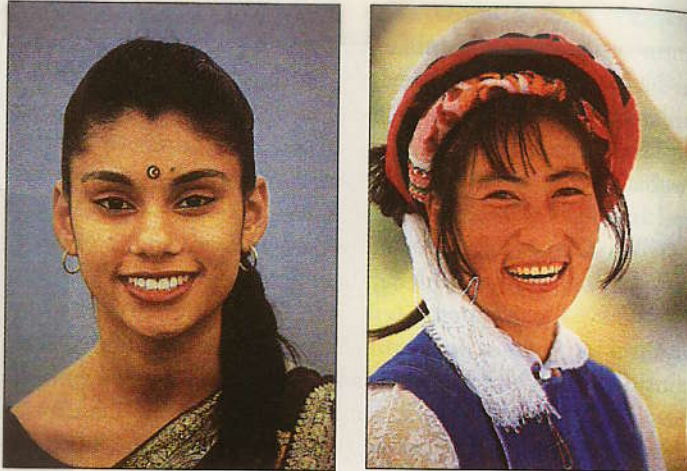


Figure 14-4 People from different geographic regions vary in such traits as height, blood type, skin color, and facial features. A. An Indian woman. B. An ethnic minority woman from China. (A, Superstock; B, Paul Greblunas/Tony Stone Images)

What Are Human "Races"?

We all know that the human species varies in appearance and in physiology geographically. People of Indian descent, for example, are recognizably different from those of Chinese descent. We can look at the shape of the nose, the eyes, and the ears, for example, and we see differences (Figure 14-4). We can look at the distribution of individual alleles and see differences. Fifty-one percent of Nigerians have type O blood, compared to only 30 percent of Japanese. Twenty percent of Russians have type B blood, while the Amerindians of Lima, Peru, have no detectable levels of type B blood at all.

Nineteenth-century anthropologists struggled to classify human groups into a few major races. Some systems identified only 12 races, while other systems listed 30 or more. One problem was that no matter how anthropologists classified humans, there always seemed to be tribes or nations that would not fit into any known group. The Basques, who live in the Pyrenees mountains between France and Spain, for example, appear European. Yet their language and culture are unlike any other in the world, and some researchers once argued that they are direct descendants of Stone Age Europeans. Similarly, the Bushmen are unique among African groups in both appearance and physiology (Figure 14-5).

A more serious problem with the grouping of humans into races is that most groups do *not* stand out from those around them; they blend. Because groups of humans inevitably mix, through migration, warfare, and trade, human "races" are never pure. Both the Japanese and British, for example, take pride

in the purity of their island races. Yet the Japanese are a grade mixture of Korean and Ainu north islanders (a people of possibly European descent). This mix shows up in the distribution of blood types from one end of Japan to the other (Figure 14-6).

The British are even more of a melting pot than the Japanese. The Bronze Age Beaker Folk mixed with the Indo-European Celts in the first thousand years B.C. In the next thousand years, the Angles, the Saxons, the Jutes, and the Picts arrived, followed by the Vikings and their descendants, the



Figure 14-5 Classifying humans into a few discrete races has been unsuccessful. Named races frequently include markedly distinct peoples. A. A Bushman from Namibia. B. A rubber plantation foreman from the Ivory Coast. (A, M.P. Kahl/Photo Researchers, Inc.; B, Charles O. Cecil/Visuals Unlimited)

SIMILARITY AND DIFFERENCES in Blood Group Allele Frequencies by "Race" or Population

How does Genetic Variation (different alleles)
VARY between & within Populations?

Examples of extreme differentiation and close similarity in blood group allele frequencies in three racial groups

Gene	Alleles	Caucasoid	Negroid	Mongoloid
Duffy	Fy	.0300	.9393	.0985
	Fy ^a	.4208	.0607	.9015
	Fy ^b	.5492	—	—
Rhesus	R ₀	.0186	.7395	.0409
	R ₁	.4036	.0256	.7591
	R ₂	.1670	.0427	.1951
	r	.3820	.1184	.0049
	r'	.0049	.0707	0
	others	.0239	.0021	0
P	P ₁	.5161	.8911	.1677
	P ₂	.4839	.1089	.8323
Auberger	Au ^a	.6213	.6419	
	Au	.3787	.3581	
Xg	Xg ^a	.67	.55	.54
	Xg	.33	.45	.46
Secretor	Se	.5233	.5727	
	se	.4767	.4273	

① Different
variation
but most
Alleles
Present

② SAME
VARIATION

Source: R. C. Lewontin, *The Genetic Basis of Evolutionary Change* (Columbia University Press, 1974).

- ① Most alleles in ALL "races" / populations
- ② NO HOMOZYGOUSITY at any locus
- ③ Some Allelic Frequencies differ between "Races" & some are the same! Duffy differs & X₂ is the same. why? Adaptive value.
- ④ Auberger, Xg, & Secretor Loci show how alleles vary within populations similarly & show NO between population differences.
- ⑤ Wide range of different Alleles within/between "Races"

SIMILARITY & Differences
in Blood Group Gene
ALLELES

VARIATION
Between Groups

Wide
VARIATION [

Wide
VARIATION [

Little
VARIATION [

26-2 TABLE Examples of Extreme Differentiation and Close Similarity in Blood Group Allelic Frequencies in Three Racial Groups

Gene	Allele	POPULATION		
		Caucasoid	Negroid	Mongoloid
Duffy	Fy	0.0300	0.9393	0.0985
	Fy ^a	0.4208	0.0000	0.9015
	Fy ^b	0.5492	0.0607	0.0000
Rhesus	R ₀	0.0186	0.7395	0.0409
	R ₁	0.4036	0.0256	0.7591
	R ₂	0.1670	0.0427	0.1951
	r	0.3820	0.1184	0.0049
	r'	0.0049	0.0707	0.0000
	Others	0.0239	0.0021	0.0000
P	P ₁	0.5161	0.8911	0.1677
	P ₂	0.4839	0.1089	0.8323
Auberger	Au ^a	0.6213	0.6419	No data
	Au	0.3787	0.3581	No data
Xg	Xg ^a	0.67	0.55	0.54
	Xg	0.33	0.45	0.46
Secretor	Se	0.5233	0.5727	No data
	se	0.4767	0.4273	No data

Source: A summary is provided in L. L. Cavalli-Sforza and W. F. Bodmer, *The Genetics of Human Populations* (W. H. Freeman and Company, 1971), pp. 724-731. See L. L. Cavalli-Sforza, P. Menozzi, and A. Piazza, *The History and Geography of Human Genes* (Princeton University Press, 1994), for detailed data.

We all have the SAME Genes
But MAY differ in Frequencies
of some alleles/genes/vars

Differences in Allelic Frequencies
in Various Populations

24-5 TABLE Frequencies of the Alleles I^A , I^B , and i at the ABO Blood Group Locus in Various Human Populations

Population	I^A	I^B	i
Eskimo	0.333	0.026	0.641
Sioux	0.035	0.010	0.955
Belgian	0.257	0.058	0.684
Japanese	0.279	0.172	0.549
Pygmy	0.227	0.219	0.554

Source: W. C. Boyd, *Genetics and the Races of Man*. D. C. Heath, 1950.

24-1 TABLE Frequencies of Genotypes for Alleles at MN Blood Group Locus in Various Human Populations

Population	GENOTYPE			ALLELE FREQUENCIES	
	M/M	M/N	N/N	$p(M)$	$q(N)$
Eskimo	0.835	0.156	0.009	0.913	0.087
Australian aborigine	0.024	0.304	0.672	0.176	0.824
Egyptian	0.278	0.489	0.233	0.523	0.477
German	0.297	0.507	0.196	0.550	0.450
Chinese	0.332	0.486	0.182	0.575	0.425
Nigerian	0.301	0.495	0.204	0.548	0.452

Source: W. C. Boyd, *Genetics and the Races of Man*. D. C. Heath, 1950.

24-2 TABLE Frequencies of Gametic Types for MNS System in Various Human Populations

Population	GAMETIC TYPE				HETEROZYGOSITY (H)	
	$M S$	$M s$	$N S$	$N s$	From gametes	From alleles
Ainu	0.024	0.381	0.247	0.348	0.672	0.438
Ugandan	0.134	0.357	0.071	0.438	0.658	0.412
Pakistani	0.177	0.405	0.127	0.291	0.704	0.455
English	0.247	0.283	0.080	0.290	0.700	0.469
Navaho	0.185	0.702	0.062	0.051	0.467	0.286

Source: A. E. Mourant, *The Distribution of the Human Blood Groups*. Blackwell Scientific, 1954.

Reasons FOR ALLELIC variation Between Populations

① Founder Effect / Geographical Isolation

↳ Selective mating do to geography/culture
or both

② Adaptive Value

↳ Have positive effect in specific
environments

e.g. Hb^S Sickle-Cell Globin allele
(India, Africa, Mediterranean)

Duffy

Skin Color Genes

GEOGRAPHY &/or ADAPTIVE VALUE

But do not vary across whole genome ->
most loci are "neutral"

NOVEL EXAMPLE OF ALLELE ADAPTIVE VALUE

- ① FyB^{ES} (Erythroid Silent) is a Major Duffy allele in African Americans & Black populations of African heritage. Rarely in other populations!
- ② Encodes chemokine receptor protein on blood cell membrane. Parasite receptors for Malarial *PLASMODIUM* bind to this receptor.
- ③ FyB^{ES} is a Switch mutation allele - FyB^{ES} allele CANNOT be transcribed \therefore no FyB or *PLASMODIUM* receptor protein -
- ④ No binding of Malaria *PLASMODIUM* - \therefore Malarial Resistance
- ⑤ High Frequency in Geographical Regions with high levels of Malaria!
- ⑥ Other Duffy alleles (Chromosome 1)
@ different frequencies in other populations

There is a Large variation in
 D1S80 VNTR alleles within
 populations but little between

Table 15.2 Allele frequencies for D1S80 among U.S. population groups

Repeat number	Caucasian	Hispanic	African American	Asian
14	0	0	0	0
15	0	0.001	0	0
16	0.001	0.010	0.002	0.034
17	0.002	0.009	0.028	0.025
18	0.237	0.224	0.073	0.152
19	0.003	0.005	0.003	0.022
20	0.018	0.013	0.032	0.007
21	0.021	0.028	0.115	0.034
22	0.038	0.024	0.081	0.017
23	0.012	0.009	0.014	0.017
24	0.378	0.315	0.234	0.230
25	0.046	0.072	0.045	0.027
26	0.020	0.007	0.006	0
27	0.007	0.016	0.008	0.047
28	0.063	0.078	0.130	0.076
29	0.052	0.055	0.053	0.042
30	0.008	0.039	0.009	0.123
31	0.072	0.053	0.054	0.093
32	0.006	0.005	0.007	0.012
33	0.003	0.004	0.004	0.005
34	0.001	0.006	0.086	0.005
35	0.003	0	0.002	0.005
36	0.004	0.011	0.001	0.005
37	0.001	0.004	0	0.007
38	0	0	0	0
39	0.003	0.004	0.003	0.005
40	0	0	0	0
41	0	0.002	0.002	0.007
>41	0.001	0.006	0.007	0.002
Sample size	718	409	606	204

Source: Data from B. Budowle, et al. 1995. *Journal of Forensic Science* 40:38

JUST
 AS
 Predicted
 FROM
 CLASS
 Genotype

VNTR
 used
 for
 HCTDA DNA
 Fingerprint

NO ADAPTIVE VALUE!!
 Good locus for Forensics -
 Group Neutral!

70

There is More Genetic Diversity
within Populations than *Between*
 Populations!! So much for the
 concept of racial "parity"!!!

Proportion of genetic diversity accounted
 for within and between populations and
 races

Gene	Total H_{species}	Proportion		
		<i>within any population</i> Within Populations	Within Races between Populations	Between Races
Hp	.994	.893	.051	.056
Ag	.994	.834	—	—
Lp	.639	.939	—	—
Xm	.869	.997	—	—
Ap	.989	.927	.062	.011
6PGD	.327	.875	.058	.067
PGM	.758	.942	.033	.025
Ak	.184	.848	.021	.131
Kidd	.977	.741	.211	.048
Duffy	.938	.636	.105	.259
Lewis	.994	.966	.032	.002
Kell	.189	.901	.073	.026
Lutheran	.153	.694	.214	.092
P	1.000	.949	.029	.022
MNS	1.746	.911	.041	.048
Rh	1.900	.674	.073	.253
ABO	1.241	.907	.063	.030
Mean		<u>.854</u>	<u>.083</u>	<u>.063</u>

*More
 genetic
 diversity
 within any
 population
 than
 between
 populations!*

Source: R. C. Lewontin, *Genetic Basis of Evolutionary Change* (Columbia University Press, 1974).

- ① 85% of Human Genetic Variation occurs *within* Populations + Between Individuals in that population!
- ② Remaining 15% of Human Genetic Variation Split Between Different Populations of Same "race" (8%) + Between Different "Races" (6%).
- ③ Only 6% of Human Genetic Variation due to Differences between Races!! *Geographic*

The SAME ALLELES are present
in most Populations but
at Different Frequencies

Mapping Human History

Mary-Claire King and Arno G. Motulsky

The DNA of modern humans contains a record of the travels and encounters of our ancestors. The genotypes of people living today are the result of ancient human migrations, the continuous appearance of new mutations, selection by climate and infection for genetic alleles that conferred a survival advantage, and mating patterns determined by cultural norms. By sampling genotypes from people across the globe, geneticists have reconstructed the major features of our history: our ancient African origin, migrations out of Africa, movements and settlements throughout Eurasia and Oceania, and peopling of the Americas (1–5). As genomic technology has improved, these analyses of genotype have successively incorporated new markers: blood groups (2), protein polymorphisms (2), mitochondrial DNA sequences (1), Y chromosome haplotypes (3), and highly variable nuclear microsatellite markers (4, 5).

The most recent contribution to this literature is by Rosenberg *et al.* (6) on page 2381 of this issue. These investigators explored the genetic structure of human pop-

ulations using highly variable markers on the human autosomes of individuals from different parts of the world. The genotyped markers were microsatellite short tandem repeat sequences that do not encode any expressed genes and are generally selectively neutral. The populations studied were defined by geography, language, and culture, and participating individuals were well rooted in their populations, with several generations of ancestors known to have lived in the same locale as the participant. Genotypes from more than a thousand individuals were evaluated by a statistical method that defines clusters of people on the basis of genetic similarity at multiple loci, without using prior information about ancestry. In this method, individuals are assigned to clusters probabilistically (5, 7). Individuals may have significant probabilities of membership in more than one cluster due either to genetic similarities of groups or to ancestral intergroup matings. The world map (see the figure) illustrates variation at one microsatellite marker in 12 populations. This marker has four common alleles, each of which appears in all populations. Rare alleles are shared by fewer populations. Few alleles are unique to only one population. No allele is population specific.

Previous genetic analyses of human

history have consistently suggested that most human genetic variation is due to differences among individuals within populations rather than to differences among populations (4, 8). The Rosenberg *et al.* analysis of many more markers and many more people confirms this result: 93 to 95% of genetic variation is due to genetic differences among individuals who are members of the same population and only 3 to 5% of genetic variation is due to differences among the major population groups.

The power of the method lies in the construction of clusters on the basis of accumulated small differences in allele frequencies across many markers and many people. Statistical clustering of genotypes—composed of 4682 alleles from 377 markers in 1056 individuals from 52 populations—yields groups corresponding to major geographic regions of the world [see figure 1 in (6)]. Creation of two clusters reflects ancient human origins in Africa and rapid expansion throughout Eurasia, and migrations to the Americas from East Asia. Creation of five clusters yields groups corresponding to five major geographic regions of the world: Africa, Eurasia (Europe, the Middle East, Central and South Asia), East Asia, Oceania, and America. There is excellent agreement between membership of individuals in these clusters and their self-identified regions of origin. Similar results were obtained by the same statistical approach based on fewer populations and fewer markers [Table 2 of (5)].

The authors are in the Department of Medicine (Medical Genetics) and Department of Genome Sciences, University of Washington, Seattle, WA 98195, USA. E-mail: mcking@u.washington.edu

We Learn about our origins
& diversity from the Human
Genome Project

SCIENCE'S COMPASS

Population substructure could be consistently identified within some geographic regions but not others. Within Africa, for example, analysis consistently yielded the same four subclusters: Mbuti Pygmies, Biaka Pygmies, San peoples, and speakers of Niger-Kordofanian languages (Bantu, Yoruba, and Mandenka populations). In contrast, within Europe, multiple analyses were not consistent. Many more individuals will need to be included to sort out European demographic history.

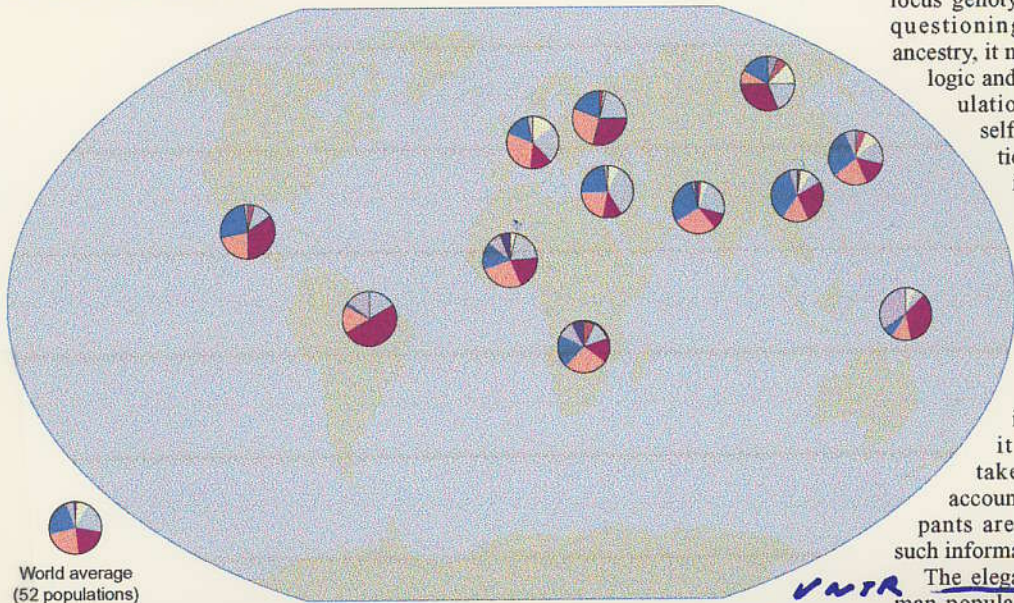
without discontinuities between clusters. After thousands of years—if enough markers and people are studied—allele frequency differences are collectively adequate to create clusters that correspond to the major migrations of human history.

What are the implications of the Rosenberg *et al.* findings for medicine? The current medical literature increasingly includes studies exploring population differences in disease incidence or in efficacy or adverse responses to drug treatment (5,

to dictate medical management for the entire group. Instead, critical alleles influencing disease risk or response to treatment are likely to be either ancient, worldwide, and relatively common in many populations, or geographically localized and individually rare (10).

Differences among populations in disease frequency and treatment outcome certainly occur but may not be genetic in origin. Given that the major population origin of groups can be defined by multi-locus genotype clustering (5–7) without questioning individuals about their ancestry, it may be tempting in epidemiologic and clinical studies to omit population characterization through self-reporting. However, correlations identified by the clustering method may be falsely ascribed to genes when in fact they had nothing to do with genetics but were caused by social, economic, or discriminatory factors limited to a genetically defined population cluster (11). To evaluate medically important group differences, it is therefore necessary to take all such risk factors into account. Patients and study participants are usually the best source of such information.

The elegant statistical analysis of human population structure by Rosenberg and colleagues reflects the major human migrations out of Africa, into Europe, across Asia, into Oceania, and to the Americas. By genotyping a large sample of an individual's alleles, it is possible to identify the migrations in which his or her ancestors participated. But the link between historical genetic demography and medically important risk is complex. Disease susceptibility may be genetic but not geographically clustered, or geographically clustered but not genetic, or neither, or both.



World average
(52 populations)

Humans on the move. Worldwide genetic variation at a neutral marker. Allele frequencies of one randomly chosen microsatellite marker reveal common alleles shared in all populations and the gradual and arbitrary differences in allele frequencies across geographic regions. Populations shown in this example are Yoruba and Bantu (Africa); French, Russians, Palestinians, and Pakistani Brahui (Eurasia); Han Chinese, Japanese, and Yakut (East Asia); New Guineans (Oceania); and Maya and Karitianans (America). Each color on the diagrams represents one of nine alleles of GGAA29H03 (D13S1493), which range in length from 219 to 255 base pairs. By accumulating small differences in allele frequencies from hundreds of such highly variable markers and hundreds of people, statistical methods reveal genetic clusters of Africans, Eurasians, East Asians, Pacific Islanders, and Americans, corresponding to major ancient human migrations (6).

The identification of clusters corresponding to the major geographic regions may depend on the sampling of individuals from well-defined, relatively homogeneous populations. If individuals were sampled from a worldwide "grid" (or a worldwide grid weighted by population density), the clusters might be much less precisely defined. Does the correspondence of worldwide genetic clusters and major geographic regions suggest borders around genetic clusters analogous to the physical borders—oceans, mountain ranges, and deserts—separating geographic regions? No. Both the results of Rosenberg and colleagues and those of previous studies (1–5, 8) indicate that unlike separations between geographic regions, differences in allele frequencies are gradual,

9). The rationale of these studies is that alleles influencing disease susceptibility or treatment response may differ in frequency across populations. Consequently, individuals would be better served if critical genotypes were taken into account when assessing disease risk or designing treatment regimens. In the absence of knowing the identities of the critical alleles, personal ancestry as indicated by study participants is often used as an initial but potentially misleading substitute.

The Rosenberg *et al.* data suggest that with the exception of ancient highly selected loci (for example, the Duffy null blood group, which confers complete protection against vivax malaria), very few alleles will be both confined to one population and common enough in that population

References

1. A. C. Wilson, R. L. Cann, *Sci. Am.* 266, 68 (April 1992).
2. L. L. Cavalli-Sforza, P. Menozzi, A. Piazza, *The History and Geography of Human Genes* (Princeton Univ. Press, Princeton, NJ, 1994).
3. P. A. Underhill *et al.*, *Nature Genet.* 26, 358 (2000).
4. G. Barbujani, A. Magagni, E. Minch, L. L. Cavalli-Sforza, *Proc. Natl. Acad. Sci. U.S.A.* 94, 4516 (1997).
5. J. F. Wilson *et al.*, *Nature Genet.* 29, 265 (2001).
6. N. A. Rosenberg *et al.*, *Science* 298, 2381 (2002).
7. J. K. Pritchard, M. Stephens, P. Donnelly, *Genetics* 155, 945 (2000).
8. R. C. Lewontin, *Evol. Biol.* 6, 381 (1972).
9. H. G. Xie, R. B. Kim, A. J. Wood, C. M. Stein, *Annu. Rev. Pharmacol. Toxicol.* 41, 815 (2001).
10. J. K. Pritchard, N. J. Cox, *Hum. Mol. Genet.* 11, 2417 (2002).
11. N. Risch, E. Burchard, E. Ziv, H. Tang, *Genome Biol.* 3, comment2007 (2002).

CONCLUSION

① If 85% of Human Genetic Variation occurs between different people within any given population (localized)

② If only 7% of Human Genetic Variation occurs between "RACES" (Novel Alleles Specific to a "race") → e.g., Fy^bE^s

③ Then Losing all "races" except one retains 94% of all Human Genetic Variation!

$$[85\% + (15\% - 7\%)] = 94\%$$

VARIAION That occurred in Ancestral Populations

85% within population genetic variability

8% between populations of SAME "race"

7% between "race" genetic variability

④ ∴ Humans highly Heterozygous or Hybrids —
* if above not true — most of us would not be here — need genetic variation to survive!

So what is a "Race"?

- ① Primarily a sociological concept — but could be a localized or inbred population that has a higher frequency of alleles at a very small number of loci. Affects few physical features.
- ② High frequency alleles in one "race" are present at lower frequency in other "races." ALL humans have SAME genes — differ in form Mostly within populations!
- ③ Heterozygosity (variation) high in human populations — ALL populations. NONE homozygous at all loci!
- ④ NO such thing as a "pure" race — would have little variation —
- ⑤ Genes affecting physical features NOT representative of genes across genome —

Geographical Ancestry is Relevant — Many "racial" groups now have multiple ancestries because of admixture & migration

∴ "Race" Classification is **Arbitrary**
& Based on a few Traits — Not
Science Based

CAN DEFINE BY MANY CRITERIA....

RACE BY RESISTANCE

Traditionally we divide ourselves into races by the twin criteria of geographic location and visible physical characteristics. But we could make an equally reasonable and arbitrary division by the presence or absence of a gene, such as the sickle-cell gene, that confers resistance to malaria.

By this reckoning, we'd place Yemenites, Greeks, New Guineans, Thai, and Dinkas in one "race," Norwegians and several black African peoples in another.

People who possess anti-malarial genes (alleles!) & those that do not

RACE BY DIGESTION

We could define a race by any geographically variable trait—for example, the retention in adulthood of the enzyme lactase, which allows us to digest milk. Using this as our divisive criterion, we can place northern and central Europeans with Arabians and such West African peoples as the Fulani; in a "lactase-negative race," we can group most other African blacks with east Asians, American Indians, southern Europeans, and Australian aborigines.

People who possess Lactase & those that do not!

RACE BY FINGERPRINTS

Probably the most trivial division of humans we could manage would be based on fingerprint patterns. As it turns out, the prevalence of certain basic features varies predictably among peoples: in the "Loops" race we could group together most Europeans, black Africans, and east Asians. Among the "Whorls" we could place Mongolians and Australian aborigines. Finally, in an "Arches" race, we could group Khoisans and some central Europeans.

People who possess Loops & those that have arches!

RACE BY GENES

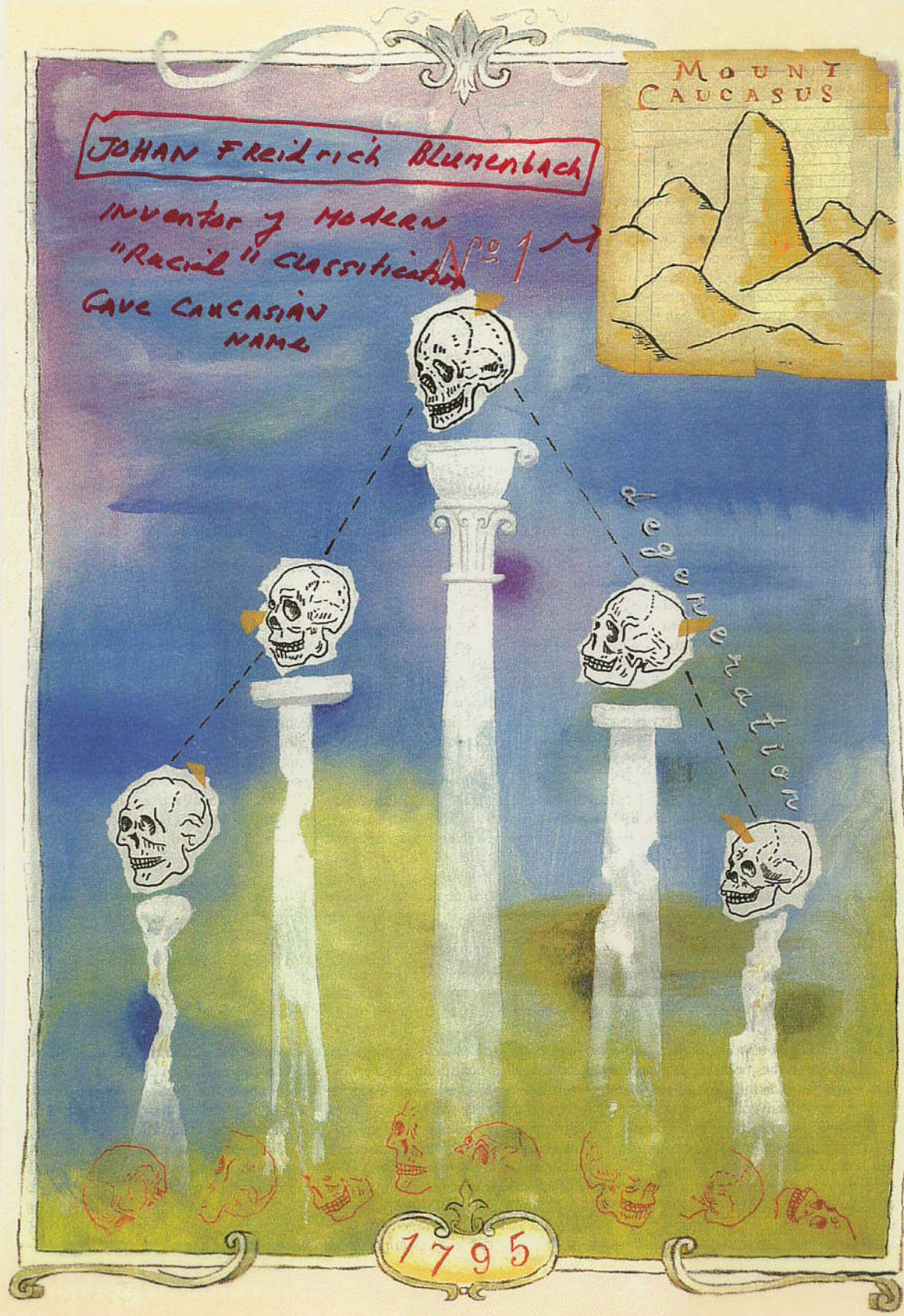
One method that seems to offer a way out of arbitrariness is to classify peoples by degree of genetic distinctness. By this standard the Khoisans of southern Africa would be in a race by themselves. African blacks would form several other distinct races. All the rest of the world's peoples—Norwegians, Navajo, Greeks, Japanese, Australian aborigines, and so on—would, despite their greatly differing external appearance, belong to a single race.

People who possess one form of an allele (localized inbred population) & those that do not

RACIAL CLASSIFICATION DIDN'T COME FROM SCIENCE BUT FROM THE BODY'S SIGNALS FOR DIFFERENTIATING ATTRACTIVE FROM UNATTRACTIVE SEX PARTNERS, AND FRIEND FROM FOE.

physical traits

HOW DID WE GET TO WHERE WE ARE?



JOHAN FRIEDRICH BLUMENBACH

- ① Changed Linnaean neutral-geography-based Human Classifications to a value-based Classification—
- ② Said all "races" originated from one place/origin - in Europe around Mt. CAUCASUS - because they are the most "beautiful" race - even though he felt all races were equal in all respects & argued with people that didn't!
- ③ Classified five "races"

CAUCASIAN
Mongolian
Ethiopian
American
Malay

Value-based Classification & putting CAUCASIANS on top - most beautiful - had disastrous consequences!

CONCLUDE

- ① Races are Arbitrary Entities - Social Constructs that are Culturally generated.
- ② Yes - there is genetic differences between "races" or relatively inbred geographical populations that can lead to physical & other differences - due to very small # genes & not reflective of whole genome
- ③ Within population genetic variation Much greater than between population genetic variation - Many loci have same allele frequencies - some differ
- ④ Only minor differences between genomes of different people or groups of people - Unity >>> Differences!
- ⑤ We are all the same - but different.
you now know why!

Race classifications Arbitrary,
unscientific, & divisive!

The End!!